

# On the Use of Large Language Models for Table Tasks

## - Introduction



Haochen Zhang

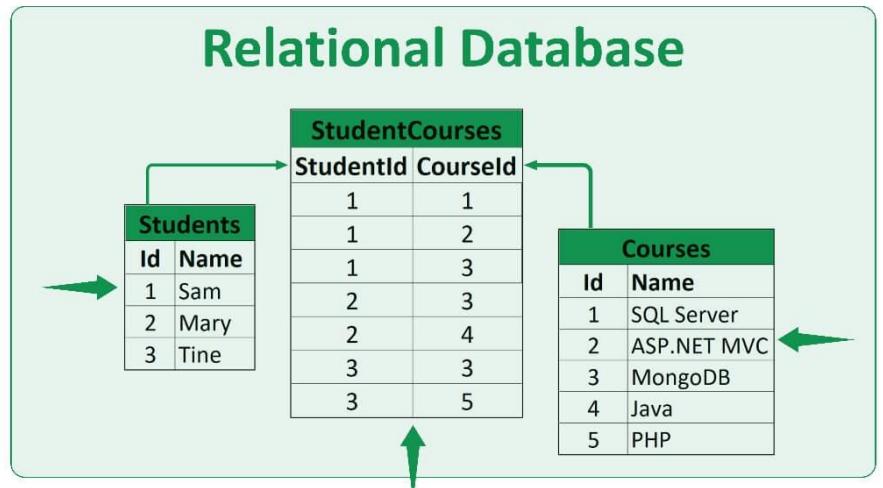


Chuan Xiao



Yuyang Dong

# Tabular data is everywhere



(a) Relational databases

"Our Data Center platform is powered by increasingly diverse drivers — demand for data processing, training and inference from large cloud-service providers and GPU-specialized ones, as well as from enterprise software and consumer internet companies. Vertical industries — led by auto, financial services and healthcare — are now at a multibillion-dollar level.

"NVIDIA RTX, introduced less than six years ago, is now a massive PC platform for generative AI, enjoyed by 100 million gamers and creators. The year ahead will bring major new product cycles with exceptional innovations to help propel our industry forward. Come join us at next month's GTC, where we and our rich ecosystem will reveal the exciting future ahead," he said.

NVIDIA will pay its next quarterly cash dividend of \$0.04 per share on March 27, 2024, to all shareholders of record on March 6, 2024.

**Q4 Fiscal 2024 Summary**

| GAAP  |          |          |         |            |             |
|---|----------|----------|---------|------------|-------------|
| (\$ in millions, except earnings per share) | Q4 FY24  | Q3 FY24  | Q4 FY23 | Q/Q        | Y/Y         |
| Revenue                                     | \$22,103 | \$18,120 | \$6,051 | Up 22%     | Up 265%     |
| Gross margin                                | 76.0%    | 74.0%    | 63.3%   | Up 2.0 pts | Up 12.7 pts |
| Operating expenses                          | \$3,176  | \$2,983  | \$2,576 | Up 6%      | Up 23%      |
| Operating income                            | \$13,615 | \$10,417 | \$1,257 | Up 31%     | Up 983%     |
| Net income                                  | \$12,285 | \$9,243  | \$1,414 | Up 33%     | Up 769%     |
| Diluted earnings per share                  | \$4.93   | \$3.71   | \$0.57  | Up 33%     | Up 765%     |

| Non-GAAP                                    |          |          |         |            |             |
|---|----------|----------|---------|------------|-------------|
| (\$ in millions, except earnings per share) | Q4 FY24  | Q3 FY24  | Q4 FY23 | Q/Q        | Y/Y         |
| Revenue                                     | \$22,103 | \$18,120 | \$6,051 | Up 22%     | Up 265%     |
| Gross margin                                | 76.7%    | 75.0%    | 66.1%   | Up 1.7 pts | Up 10.6 pts |
| Operating expenses                          | \$2,210  | \$2,026  | \$1,775 | Up 9%      | Up 25%      |
| Operating income                            | \$14,749 | \$11,557 | \$2,224 | Up 28%     | Up 563%     |
| Net income                                  | \$12,839 | \$10,020 | \$2,174 | Up 28%     | Up 491%     |
| Diluted earnings per share                  | \$5.16   | \$4.02   | \$0.88  | Up 28%     | Up 486%     |

(b) Rich documents, PDF

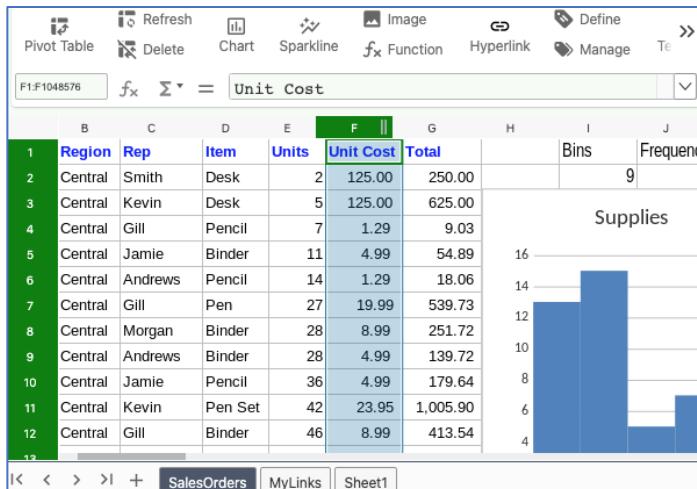
List of Large Language Models [edit]

See also: List of chatbots

For the training cost column, 1 petaFLOP-day = 1 petaFLOP/sec × 1 day = 8.64E19 FLOP. Also, only the largest model's cost is written.

| Name    | Release date <sup>[a]</sup> | Developer  | Number of parameters (billion) <sup>[b]</sup> | Corpus size   | Training cost (petaFLOP-day) | License <sup>[c]</sup>      | Notes   |
|---------|-----------------------------|------------|---|---|------------------------------|-----------------------------|---|
| GPT-1   | June 2018                   | OpenAI     | 0.117   |   | 1 <sup>[141]</sup>           | MIT <sup>[142]</sup>        | First GPT model, decoder-only transformer. Trained for 30 days on 8 P600 GPUs.  |
| BERT    | October 2018                | Google     | 0.340 <sup>[143]</sup>                        | 3.3 billion words <sup>[143]</sup>                          | 9 <sup>[144]</sup>           | Apache 2.0 <sup>[145]</sup> | An early and influential language model. <sup>[4]</sup> Encoder-only and thus not built to be prompted or generative. <sup>[146]</sup> Training took 4 days on 64 TPUv2 chips. <sup>[147]</sup>                       |
| T5      | October 2019                | Google     | 11 <sup>[148]</sup>                           | 34 billion tokens <sup>[148]</sup>                          |                              | Apache 2.0 <sup>[149]</sup> | Base model for many Google projects, such as Imagen. <sup>[150]</sup>   |
| XLNet   | June 2019                   | Google     | 0.340 <sup>[151]</sup>                        | 33 billion words  | 330                          | Apache 2.0 <sup>[152]</sup> | An alternative to BERT; designed as encoder-only. Trained on 512 TPU v3 chips for 5.5 days. <sup>[153]</sup>  |
| GPT-2   | February 2019               | OpenAI     | 1.5 <sup>[154]</sup>                          | 40GB <sup>[155]</sup> (~10 billion tokens) <sup>[156]</sup> | 28 <sup>[157]</sup>          | MIT <sup>[158]</sup>        | Trained on 32 TPUv3 chips for 1 week. <sup>[157]</sup>  |
| GPT-3   | May 2020                    | OpenAI     | 175 <sup>[50]</sup>                           | 300 billion tokens <sup>[156]</sup>                         | 3640 <sup>[159]</sup>        | proprietary                 | A fine-tuned variant of GPT-3, termed GPT-3.5, was made available to the public through a web interface called ChatGPT in 2022. <sup>[160]</sup>  |
| GPT-Neo | March 2021                  | EleutherAI | 2.7 <sup>[161]</sup>                          | 825 GiB <sup>[162]</sup>                                    |                              | MIT <sup>[163]</sup>        | The first of a series of free GPT-3 alternatives released by EleutherAI. GPT-Neo outperformed an equivalent-size GPT-3 model on some benchmarks, but was significantly worse than the largest GPT-3. <sup>[163]</sup> |

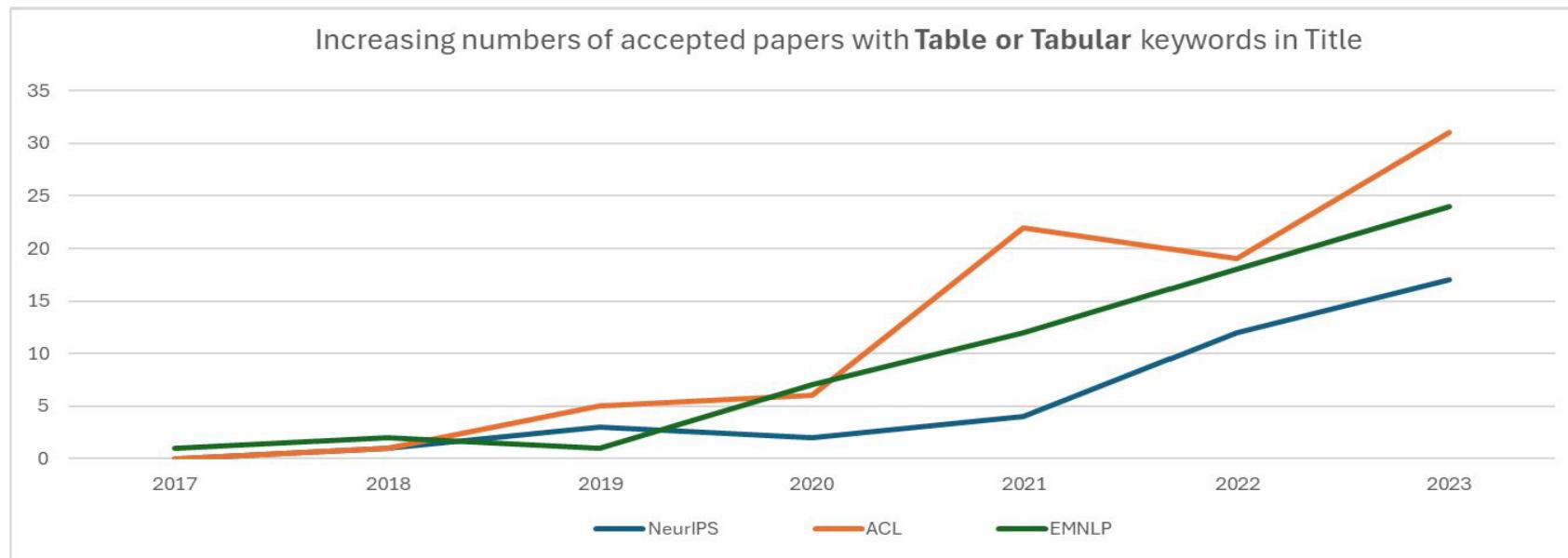
(c) webpages



(d) spreadsheet

# Growing research focus

- Grow quickly in DB, AI and NLP communities



- Recent tutorials
  - [1] Web table extraction, retrieval and augmentation, SIGIR19
  - [2] From Tables to Knowledge: Recent Advances in Table Understanding, KDD21
  - [3] Transformers for Tabular Data Representation: A tutorial on Models and Applications VLDB22, SIGMOD23
  - [4] Large Language Models for Tabular Data: Progresses and Future Directions, SIGIR24
  - A-Paper-List-of-Awesome-Tabular-LLMs, <https://github.com/SpursGoZmy/Awesome-Tabular-LLMs>

# Table tasks & benchmarks

“Prepare tables”

## Table reprocessing

- Table matching
  - Entity matching
  - Schema matching
- Table cleaning
  - Error detection
  - Data imputation
- Table augmentation
  - Row population
  - Schema augmentation
- Table search
- Table transformation

“Understand tables”

## Table understanding

- Table Interpretation
  - Entity Linking
  - Column Type Annotation
  - Relation Extraction
- Table detection

“Get answer from tables”

## Table analysis

- Table QA
- Table fact verification
- Table-to-text
- Text-to-SQL

# Table matching

“Matching two rows”

|     |  |  |  |
|-----|--|--|--|
| AAA |  |  |  |
| BBB |  |  |  |
| CCC |  |  |  |

**Entity matching**

EEE

AAA'

DDD

“Matching two columns”

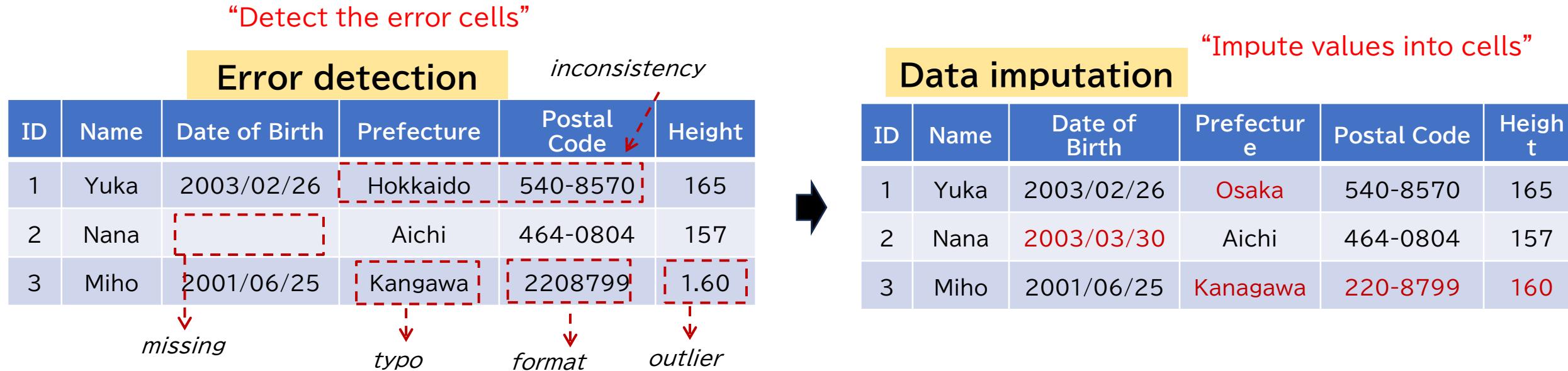
Schema  
matching

| id | name  | loc | # of employee |
|----|-------|-----|---------------|
| 1  | Apple | CA  | 154,000       |
| 2  | IBM   | NY  | 282,000       |

| id | name      | rev    |
|----|-----------|--------|
| 1  | IBM Corp  | \$57B  |
| 2  | Apple Inc | \$366B |
| 3  | GE        | \$74B  |

- Dataset and benchmark
  1. Can Foundation Models Wrangle Your Data? [VLDB23]. <https://arxiv.org/abs/2205.09911>
  2. Jellyfish: Instruction-Tuning Local Large Language Models for Data Preprocessing [EMNLP24] <https://arxiv.org/abs/2312.01678>

# Table cleaning



- Dataset and benchmark
  1. Can Foundation Models Wrangle Your Data? [VLDB23]. <https://arxiv.org/abs/2205.09911>
  2. Jellyfish: Instruction-Tuning Local Large Language Models for Data Preprocessing [EMNLP24] <https://arxiv.org/abs/2312.01678>

# Table augmentation

“Add columns/rows to table”

Column  
population

| id | name  | loc | # of employee | + |
|----|-------|-----|---------------|---|
| 1  | Apple | CA  | 154,000       |   |
| 2  | IBM   | NY  | 282,000       |   |
| +  |       |     |               |   |

Row population

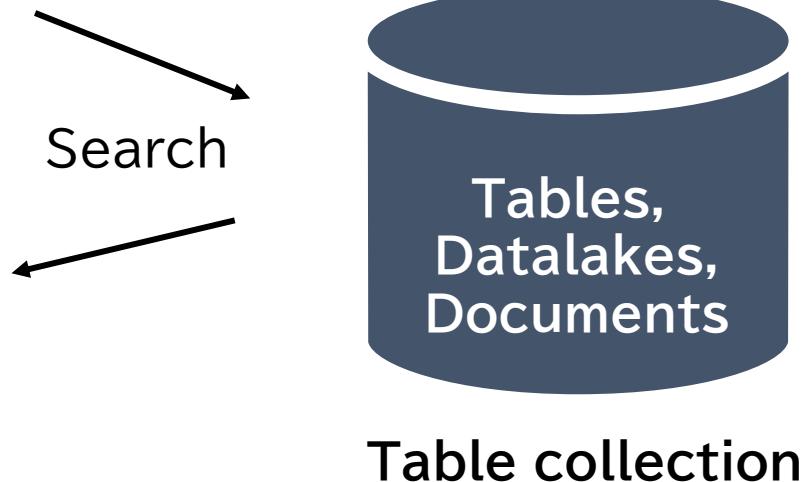
- Dataset and benchmark
1. TURL: Table Understanding through Representation Learning. [VLDB20]  
<https://arxiv.org/abs/2006.14806>
  2. TableLlama: Towards Open Large Generalist Models for Tables. [NAACL24].  
<https://osu-nlp-group.github.io/TableLlama/>

# Table search

“Retrieve tables with an NL query”

“Tables contains  
information of Apple Inc.”

| id | name  | loc | # of employee |
|----|-------|-----|---------------|
| 1  | Apple | CA  | 154,000       |
| 2  | IBM   | NY  | 282,000       |



- Dataset and benchmark
  - 1. Open-Domain Table Retrieval for Natural Questions. <https://github.com/zorazrw/nqt-retrieval>
  - 2. Open-WikiTable :Dataset for Open Domain Question Answering with Complex Reasoning over Table [EMNLP23] [https://github.com/sean0042/Open\\_WikiTable](https://github.com/sean0042/Open_WikiTable)

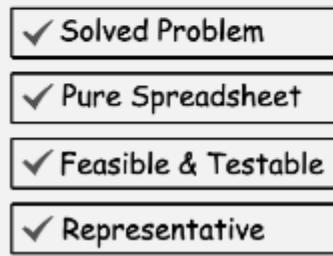
# Table transformation

“Manipulate table into wanted styles”

## 1. Data Sourcing



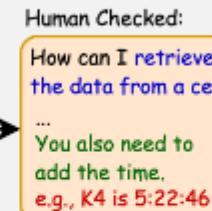
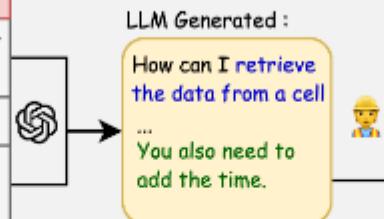
## 2. Data Filtering



## 3. Data Formatting

### Instruction Generation

| Spreadsheet Forum Post |  |
|------------------------|--|
| #1                     | I am looking for a formula that retrieve the data from a cell... |
| #2                     | Maybe: =IF(E4="", "", E4+IF...                                   |
| #3                     | Thanks, but I also need to add the time, e.g., K4 is 5:22:46     |



## 4. Testcase Construction

| A | B         | C   | D      |
|---|-----------|-----|--------|
| 1 | Name      | Age | Gender |
| 2 | Ken       | 12  | Male   |
| 3 | Bob       | 31  | Male   |
| 4 | June      | 22  | Female |
| 5 | Yang Ming | 16  | Male   |
| 6 | Jun Zhu   | 18  | Female |

apply solution:  
`=IF(B2<18, "no", "yes")`

| A | B         | C   | D      |
|---|-----------|-----|--------|
| 1 | Name      | Age | Gender |
| 2 | Ken       | 12  | Male   |
| 3 | Bob       | 31  | Male   |
| 4 | June      | 22  | Female |
| 5 | Yang Ming | 16  | Male   |
| 6 | Jun Zhu   | 18  | Female |

modify: cell B3, B5

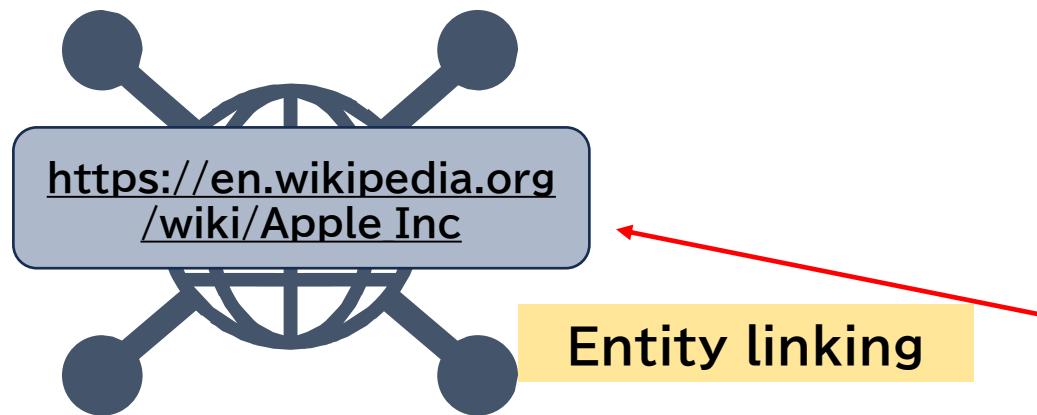
| A | B         | C   | D      |
|---|-----------|-----|--------|
| 1 | Name      | Age | Gender |
| 2 | Ken       | 12  | Male   |
| 3 | Bob       | 13  | Male   |
| 4 | June      | 22  | Female |
| 5 | Yang Ming | 16  | Male   |
| 6 | Jun Zhu   | 18  | Female |

## • Dataset and benchmark

1. SpreadsheetBench: Towards Challenging Real World Spreadsheet Manipulation [NeurIPS24] <https://spreadsheetbench.github.io/>

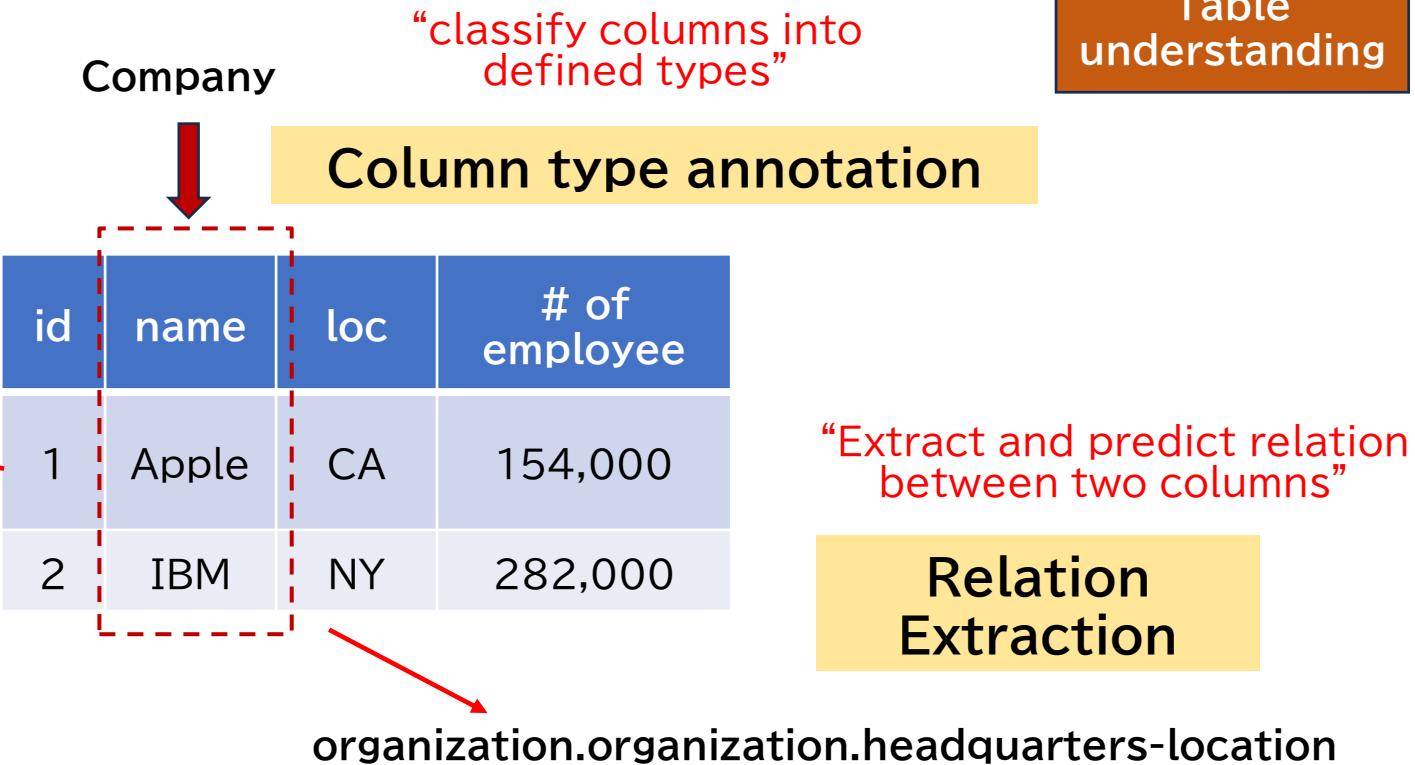
# Table Interpretation

Table understanding



Entity linking

“Match entity to knowledge base”



- Dataset and benchmark
1. TURL: Table Understanding through Representation Learning. [VLDB20]  
<https://arxiv.org/abs/2006.14806>
  2. TableLlama: Towards Open Large Generalist Models for Tables. [NAACL24].  
<https://osu-nlp-group.github.io/TableLlama/>
  3. Column Type Annotation using ChatGPT [arxiv24] <https://arxiv.org/abs/2306.00745>

# Table detection

“detect table region, structure and content”

## Table Detection

**Table 4** Multivariate analysis of factors associated with prevalence of developmental dental hard tissue anomalies ( $N = 1955$ )

|                            | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

and even lower than the caries prevalence in many other developing countries. The risk factors for caries are also more well understood [32]. This study provides evidence that the presence of developmental dental hard tissue anomalies does not increase the probability of children having caries in the study population.

Of importance is the significant association between developmental dental hard tissue anomalies and poor oral hygiene. The risk of having developmental dental hard tissue anomalies is higher in teeth cleaning [22]. It also increases molar occlusion, which also increases the risk for plaque retention [33]. The significant association between caries and this study is therefore consistent with prior observations [44, 45] and has programmatic implications for managing dental caries in children with developmental dental hard tissue anomalies. It should be treated as having high risk for poor oral hygiene and should therefore be treated more aggressively for initial visits with particular emphasis on educating the above mentioned groups about the possible use of adjunctive therapies. This is important as oral health affect adolescents perception of body image, self-esteem and social interactions [46, 47].

This study found a non-significant association between caries and presence of enamel hypoplasia under the findings of previous studies [48, 49]. However, Ferreira et al.'s [51] meta-analysis strongly indicates that developmental defects of the enamel such as enamel hypoplasia is a risk factor for caries. The presence of caries indicates that enamel hypoplasia is not a risk factor for caries in the study population from a sub-saharan developing country, while it is a risk factor in the rest of the world [52]. However, the non-significant association between developmental dental hard tissue anomalies and caries

and the significant association between developmental dental hard tissue anomalies and oral hygiene suggests a highly plausible pathophysiology for the association with developmental dental hard tissue anomalies; caries results as a secondary outcome of poor oral hygiene and not through a direct pathway. This provides evidence that caries is a secondary outcome of poor oral hygiene and further supports the hypothesis that multiple interrelated factors that may increase the susceptibility of teeth with developmental dental hard tissue anomalies to caries.

The study found gender and socioeconomic class differences in the prevalence of enamel hypoplasia differed from the findings of Babiker et al. [53] in Spain, where they found no significant difference in the prevalence of developmental dental hard tissue anomalies between males and females. This may be due to the fact that the prevalence of developmental defects of the enamel with decreasing socioeconomic status had been established, with this difference being greater for dental visits with participants receiving more education [54]. However, this study found, even though the prevalence of developmental defects of the enamel by gender remains unclear with authors identifying male at greater risk [45, 56], while others show no gender association [59, 60]. Many of these studies assessed enamel defects, regardless of whether it was primary or secondary to caries.

This study was a school based study implying that children in Southwestern Nigeria who do not attend schools are excluded from the study. This may indicate that a high proportion of children in Nigeria are not out of school [61]. Therefore the generalization of the findings to the entire population of Nigeria is limited. In addition, the study did not include the entire population of the study; the data still provides useful information highlighting the prevalence of developmental dental hard tissue anomalies.

## Table Structure Recognition

**Column**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

**Row**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

**Column Header Cell**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

**Projected Row Header Cell**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

**Text Cell**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

**Spans Cell**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

**Grid Cell**

| Variables                  | Adjusted Prevalence Ratio (APR) | Std. Err. | P value | 95 % Conf. Interval |
|----------------------------|---------------------------------|-----------|---------|---------------------|
| Oral hygiene status        |                                 |           |         |                     |
| Good oral hygiene status   | 1.00                            | -         | -       | -                   |
| Fair oral hygiene status   | 0.02                            | 0.002     | 0.14    | -0.007 - 0.05       |
| Poor oral hygiene status   | 0.07                            | 0.003     | 0.002   | 0.03 - 0.12         |
| Caries status              |                                 |           |         |                     |
| Absence of caries          | 1.00                            | -         | -       | -                   |
| Presence of caries         | 0.005                           | 0.002     | 0.77    | -0.03 - 0.04        |
| Gender                     |                                 |           |         |                     |
| Male                       | 1.00                            | -         | -       | -                   |
| Female                     | -0.006                          | 0.001     | 0.64    | -0.03 - 0.02        |
| Socioeconomic status       |                                 |           |         |                     |
| High socioeconomic class   | 1.00                            | -         | -       | -                   |
| Middle socioeconomic class | -0.001                          | 0.001     | 0.95    | -0.03 - 0.03        |
| Low socioeconomic class    | -0.007                          | 0.002     | 0.68    | -0.04 - 0.03        |

## Dataset and benchmark

1. PubTables-1M: Towards comprehensive table extraction from unstructured. [CVPR22] <https://github.com/microsoft/table-transformer>
2. TableFormer: Table Structure Understanding with Transformers. [CVPR22]. <https://github.com/IBM/TabFormer>
3. Docling Technical Report. [arxiv24] <https://github.com/DS4SD/docling-ibm-models>

# Table QA

Table analysis

“Question-answering on tabular contents”

| id | name  | loc | # of employee |
|----|-------|-----|---------------|
| 1  | Apple | CA  | 154,000       |
| 2  | IBM   | NY  | 282,000       |

Question 1: Where is the location of IBM?

Answer 1: New York

Question 2: What is the sum of employee in Apple and IBM?

Answer 2: 436,000

Table 1: An overview of table QA datasets. The representative methods without marks (e.g.  $\dagger\star\ddagger$ ) can be used on the datasets aligned in the same horizontal zone, and the methods with marks are currently adopted on the datasets with the same mark.

|                 | Dataset                    | Closed-domain | Question Type   | Representative Methods                              |
|-----------------|----------------------------|---------------|-----------------|---|
| Table -only     | WTQ $\star$ [36]           | Yes           | Factoid         | Semantic parsing-based [10,11,12,16,19,20,28,34,36] |
|                 | SQA $\ddagger\star$ [20]   | Yes           | Factoid         | [38,41,42,44,46,50]                                 |
|                 | WikiSQL $\star$ [50]       | Yes           | Factoid         | Generative method $\ddagger$ [31]                   |
|                 | Spider [47]                | Yes           | Factoid         | Matching-based method $\dagger$ [15]                |
|                 | HiTab [8]                  | Yes           | Factoid         | Extractive method $\star$ [18]                      |
|                 | AIT-QA $\dagger\star$ [23] | Yes           | Factoid         |   |
| Non-table -only | FeTaQA [32]                | Yes           | Free form       | Generative method [32]                              |
|                 | FinQA [7]                  | Yes           | Factoid         | Semantic parsing-based [7]                          |
|                 | TAT-QA [52]                | Yes           | Factoid         | Extractive methods [6,13,52]                        |
|                 | HybridQA [6]               | Yes           | Factoid         |   |
|                 | TabMCQ [22]                | Yes           | Multiple choice | Matching-based methods [22,27]                      |
|                 | GeoTSQA [27]               | Yes           | Multiple choice |   |
|                 | OTTQA [5]                  | No            | Factoid         | Retriever-reader-based methods [5,17,25,35,51]      |
|                 | NQ-tables [17]             | No            | Factoid         |   |

- Dataset and benchmark

1. TableLlama: Towards Open Large Generalist Models for Tables. [NAACL24]. <https://osu-nlp-group.github.io/TableLlama/>
2. <https://huggingface.co/datasets/SpursgoZmy/IFT-Data-For-Tabular-Tasks>

# Table fact verification

Table analysis

| id | name  | loc | # of employee |
|----|-------|-----|---------------|
| 1  | Apple | CA  | 154,000       |
| 2  | IBM   | NY  | 282,000       |

“Verify a given sentences according to table”

“IBM and Apple are U.S. companies.” → Entailed

“Apple has more employee than IBM.” → Refuted

- **Dataset and benchmark**

1. TabFact: A Large-scale Dataset for Table-based Fact Verification.[ICLR20]  
<https://tabfact.github.io/>
2. FEVEROUS: Fact Extraction and VERification Over Unstructured and Structured information. [ACL21 workshop]. <https://fever.ai/dataset/feverous.html>

# Table-to-text

Table analysis

“Generate description of table”

| id | name  | loc | # of employee |
|----|-------|-----|---------------|
| 1  | Apple | CA  | 154,000       |
| 2  | IBM   | NY  | 282,000       |



“The table presents information about two major companies, Apple and IBM, along with their locations and employee counts. Apple, headquartered in California (CA), employs 154,000 people. On the other hand, IBM, based in New York (NY), has a significantly larger workforce, with 282,000 employees.”

- Dataset and benchmark
  - 1. Neural Text Generation from Structured Data with Application to the Biography Domain [EMNLP16] <https://github.com/DavidGrangier/wikipedia-biography-dataset>
  - 2. ToTTo: A Controlled Table-To-Text Generation Dataset [EMNLP20] <https://huggingface.co/datasets/google-research-datasets/totto>
  - 3. Table-to-text: Describing table region with natural language. [AAAI18] <https://github.com/msra-nlc/Table2Text>

“Convert natural language to SQL query”

Text: How many employees in Apple?



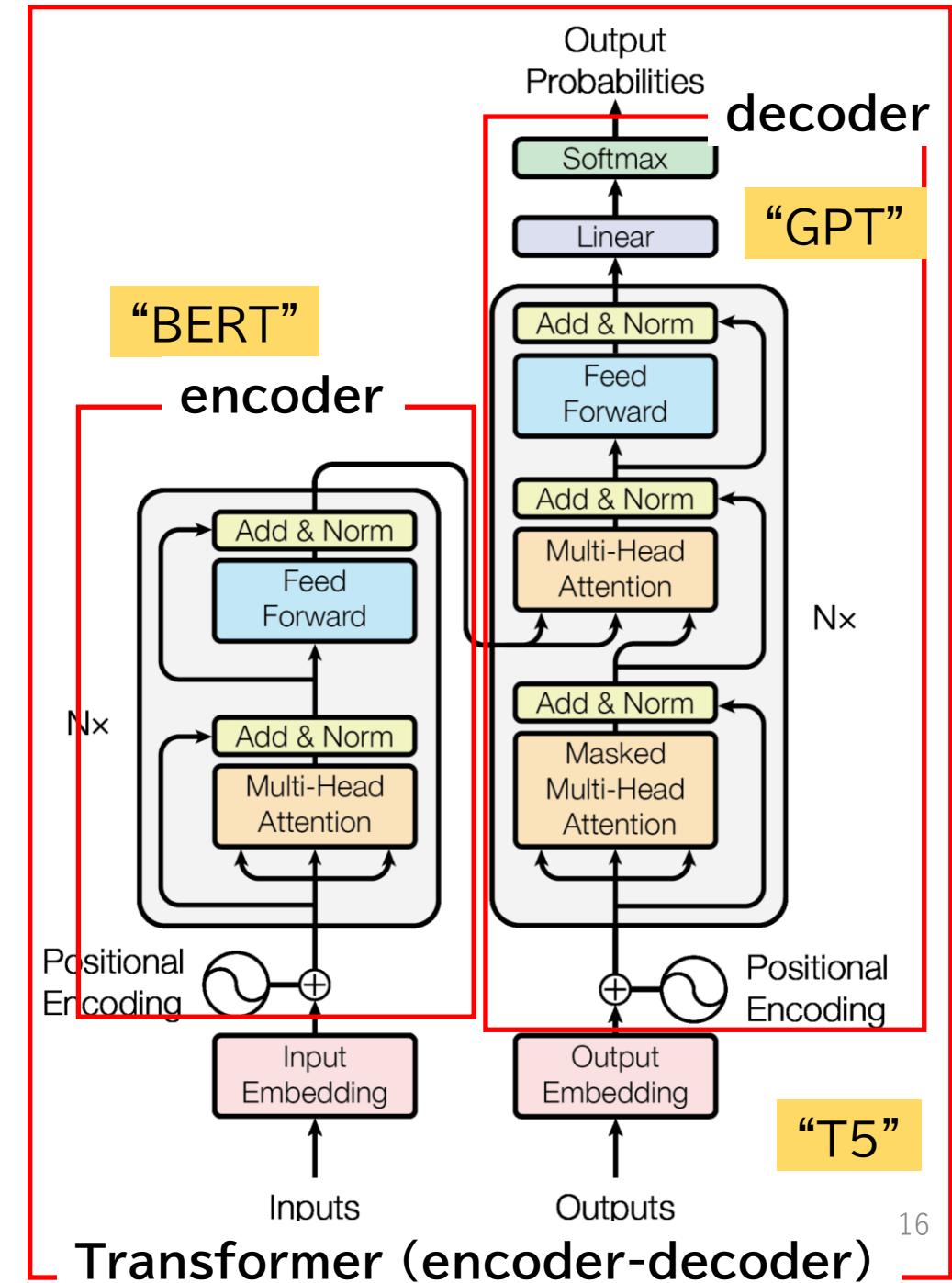
```
SELECT `employee_num`  
FROM table_name  
WHERE name = 'Apple';
```

| id | name  | loc | employee_num |
|----|-------|-----|--------------|
| 1  | Apple | CA  | 154,000      |
| 2  | IBM   | NY  | 282,000      |

- Dataset and benchmark
1. Seq2SQL: Generating Structured Queries from Natural Language using Reinforcement Learning.[ICLR18] <https://github.com/salesforce/WikiSQL>
  2. Spider [EMNLP18] <https://yale-lily.github.io/spider>

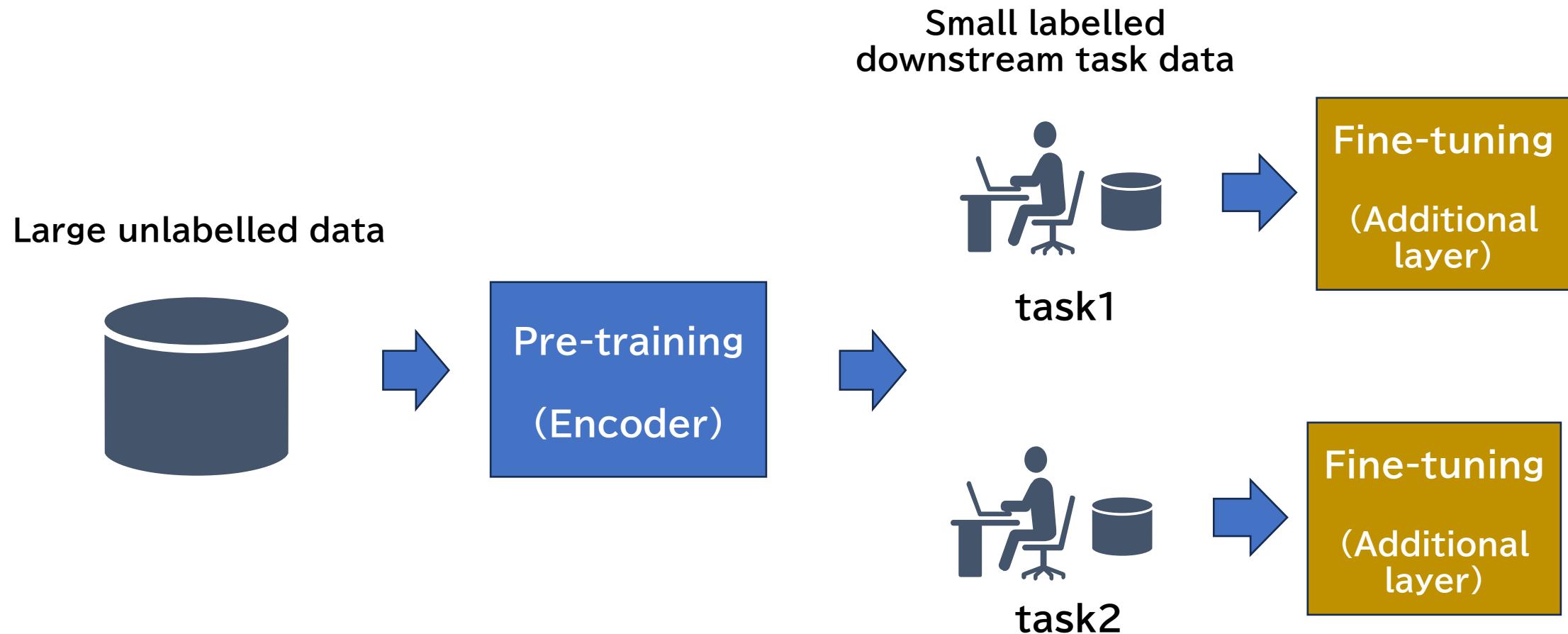
# Methods before LLM

- Rule, ML, NN-based -> skip
- Transformer-based (2018-)
  - Encoder
  - Encoder-Decoder



# Motivation of Encoder for tables

- Pretrain-and-finetune (“BERT-way”)
  - Learning good table representation (embedding) with table pretraining tasks
  - Finetune on downstream tasks



# Table-only Pretraining

- Pretrain with table contents
- TURL [VLDB20]
  - Masked Language Model (MLM)
  - Masked Entity Recovery
- TABBIE [NAACL21]
  - Detect corrupted cells
- TUTA [KDD21]
  - MLM
  - Cell filling
  - Context selection

TABBIE

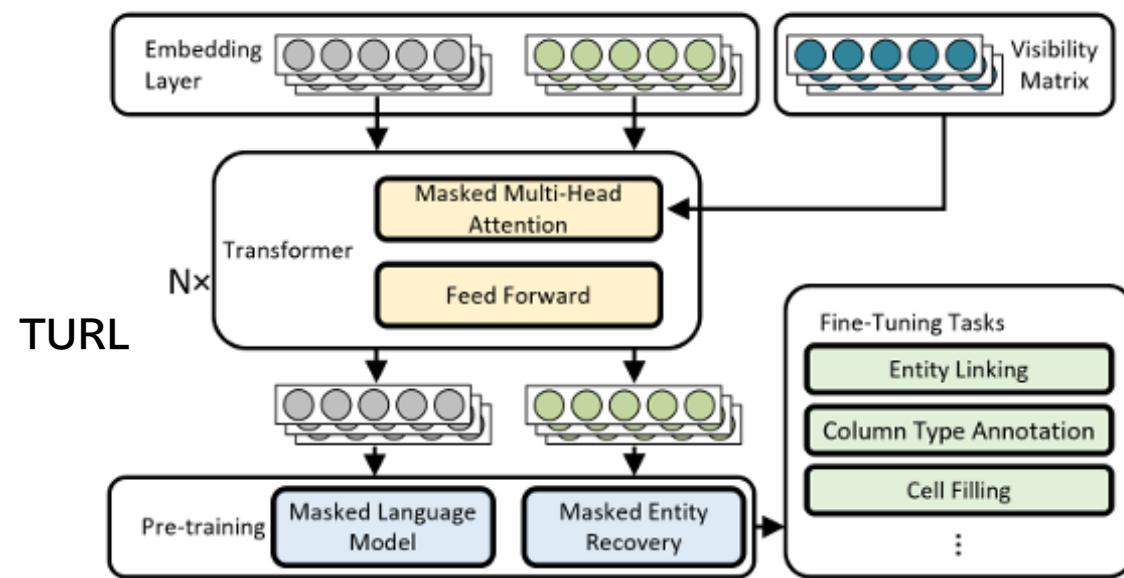
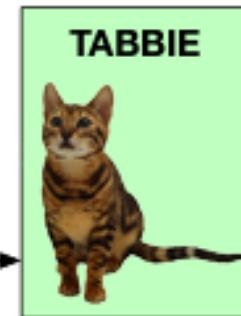


Figure 2: Overview of our TURL framework.

step 1: corrupt  
15% of cells

| Size   | Medals |
|--------|--------|
| France | 3.6    |
| Italy  | 5      |
| Spain  | 4      |



step 2: embed the table with TABBIE

step 3: train TABBIE to identify the corrupted cells

| corrupt! | real     |
|----------|----------|
| real     | corrupt! |
| real     | real     |
| real     | real     |

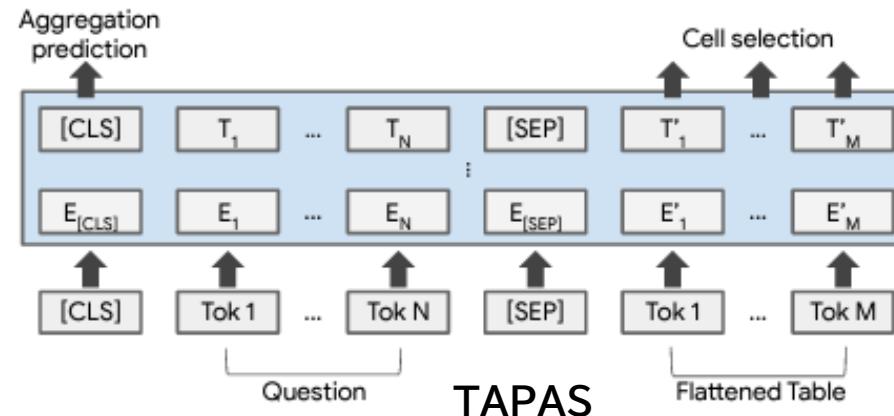
# Table-and-query Pretraining

- Pretrain with table contents & query
- TAPAS[ACL20]
  - Query and whole table
    - Aggregation prediction
    - Cell selection task
- TaBERT[NAAACL21]
  - Query and related rows
    - Masked Language Model

| op    | $P_s(op)$ | compute( $op, P_s, T$ )      |
|-------|-----------|------------------------------|
| NONE  | 0         | -                            |
| COUNT | 0.1       | .9 + .9 + .2 = 2             |
| SUM   | 0.8       | .9×37 + .9×31 + .2×15 = 64.2 |
| AVG   | 0.1       | 64.2 ÷ 2 = 32.1              |

$$S_{pred} = .1 \times 2 + .8 \times 64.2 + .1 \times 32.1 = 54.8$$

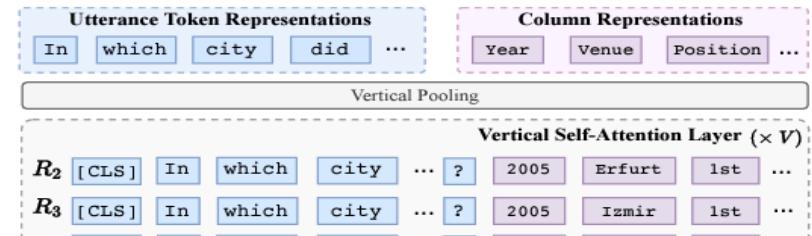
| Rank | ... | Days | $P_s$ |
|------|-----|------|-------|
| 1    | ... | 37   | 0.9   |
| 2    | ... | 31   | 0.9   |
| 3    | ... | 17   | 0     |
| 4    | ... | 15   | 0.2   |
| ...  | ... | ...  | 0     |



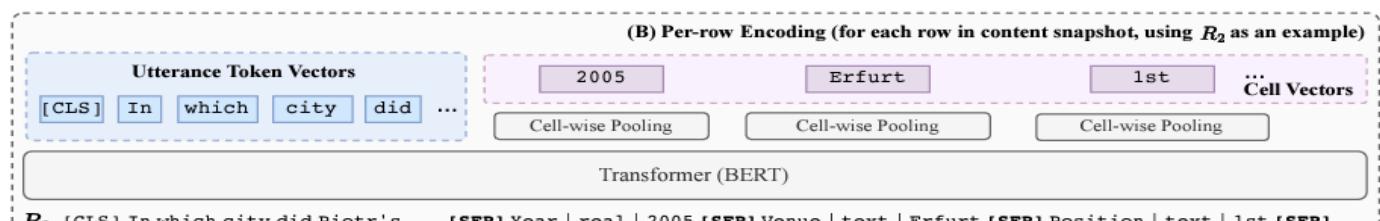
| In which city did Piotr's last 1st place finish occur? |         |          |                           |
|--|---------|----------|---------------------------|
| Year   | Venue   | Position | Event                     |
| R1 2003  | Tampere | 3rd      | EU Junior Championship    |
| R2 2005  | Erfurt  | 1st      | EU U23 Championship       |
| R3 2005  | Izmir   | 1st      | Universiade               |
| R4 2006  | Moscow  | 2nd      | World Indoor Championship |
| R5 2007  | Bangkok | 1st      | Universiade               |

Selected Rows as Content Snapshot : {R<sub>2</sub>, R<sub>3</sub>, R<sub>5</sub>}

(A) Content Snapshot from Input Table



(C) Vertical Self-Attention over Aligned Row Encodings



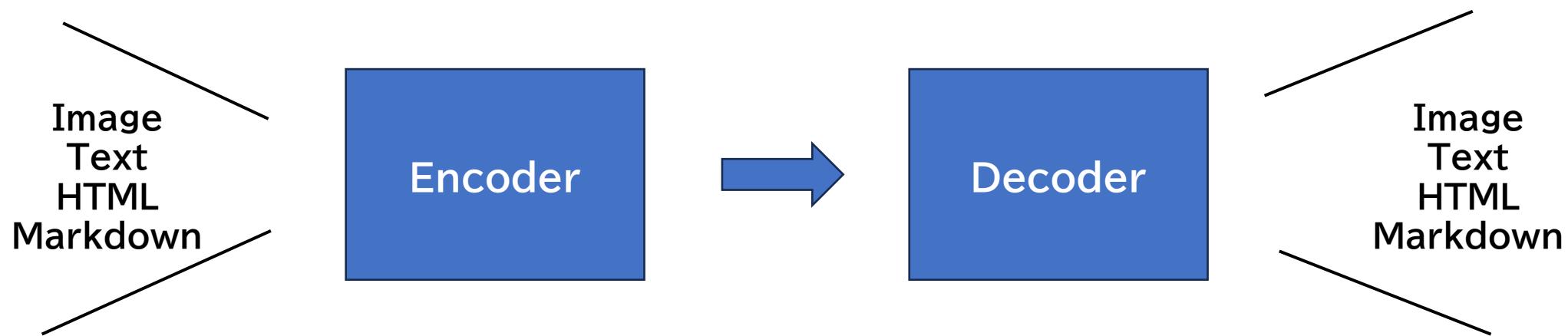
TaBERT

TAPAS: <https://arxiv.org/abs/2004.02349>

TaBERT: <https://arxiv.org/abs/2005.08314>

# Motivation of Encoder-decoder for tables

- Flexible input and output
  - Table to text
  - Text to sql
  - Table summarize
  - Table to markdown, html
- Multimodal ability
  - Image encoder -> text decoder: table OCR task, table VQA
- Generalized and good generation ability



# Text-to-text Encoder Decoder

- Generalized and Good generative ability by fine-tuning on pretrained encoder decoder model

- UnifiedSKG [EMNLP22]
  - Fine tune T5

- TaPEX [ICLR22]
  - Fine tune BART

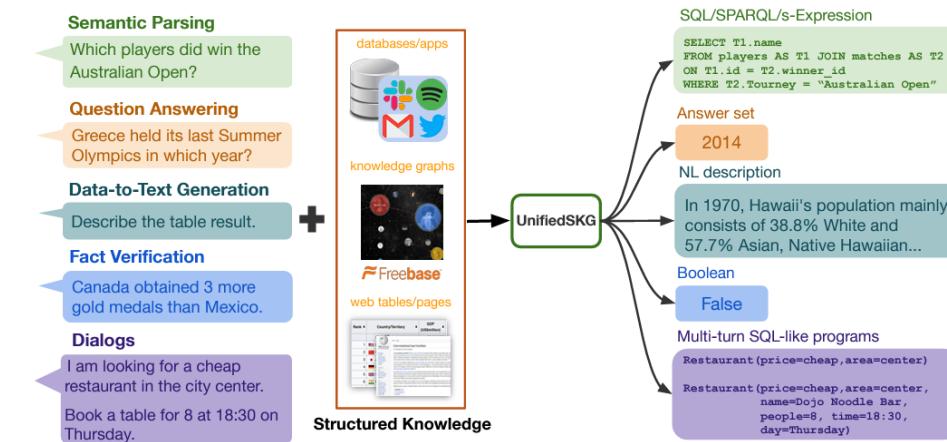
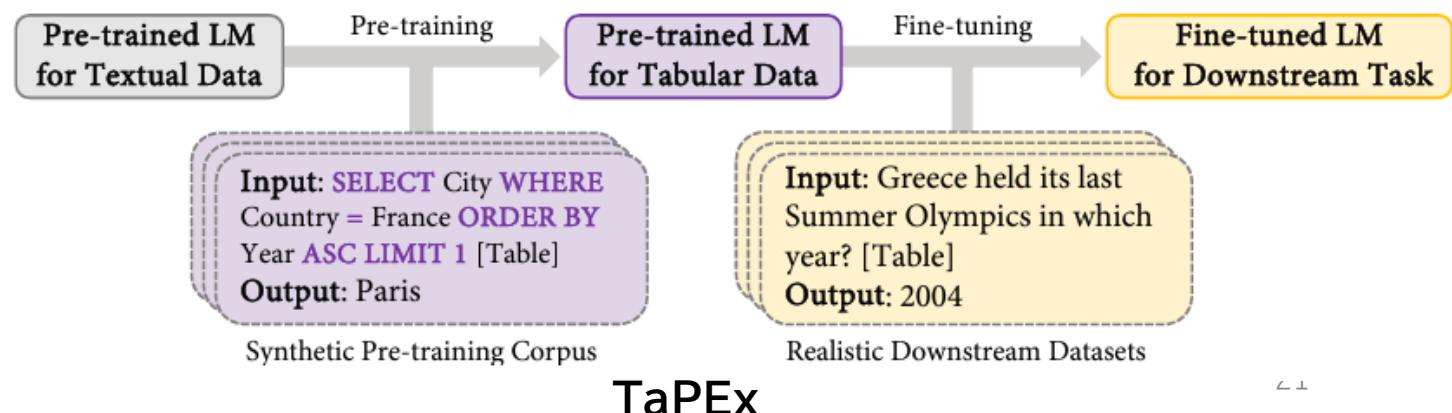


Figure 1: Structured knowledge grounding (SKG) leverages structured knowledge to complete user requests. By casting inputs and outputs into the text-to-text format, UNIFIEDSKG standardizes datasets, models, code, experiments, and metrics for 21 SKG tasks.

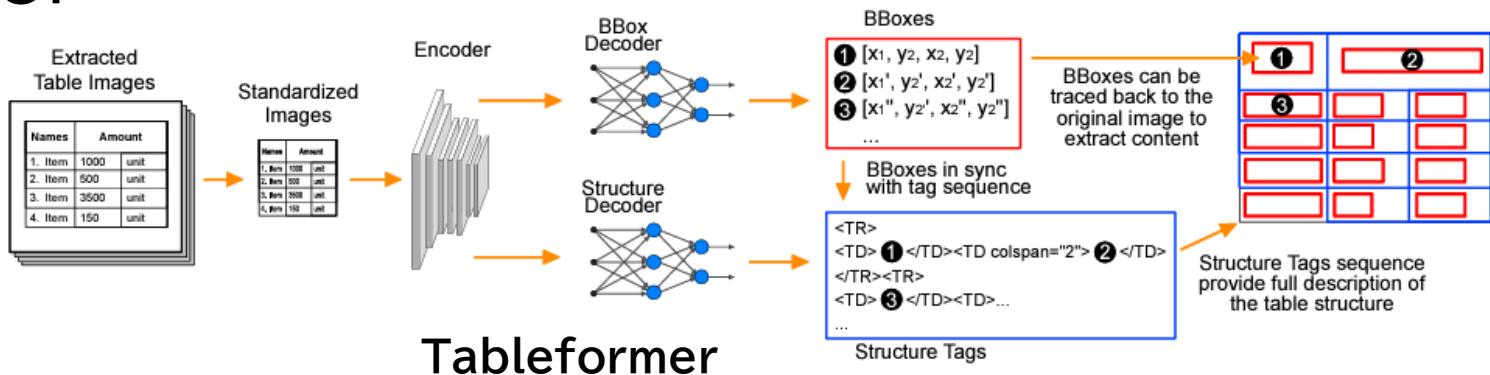
## UnifiedSKG



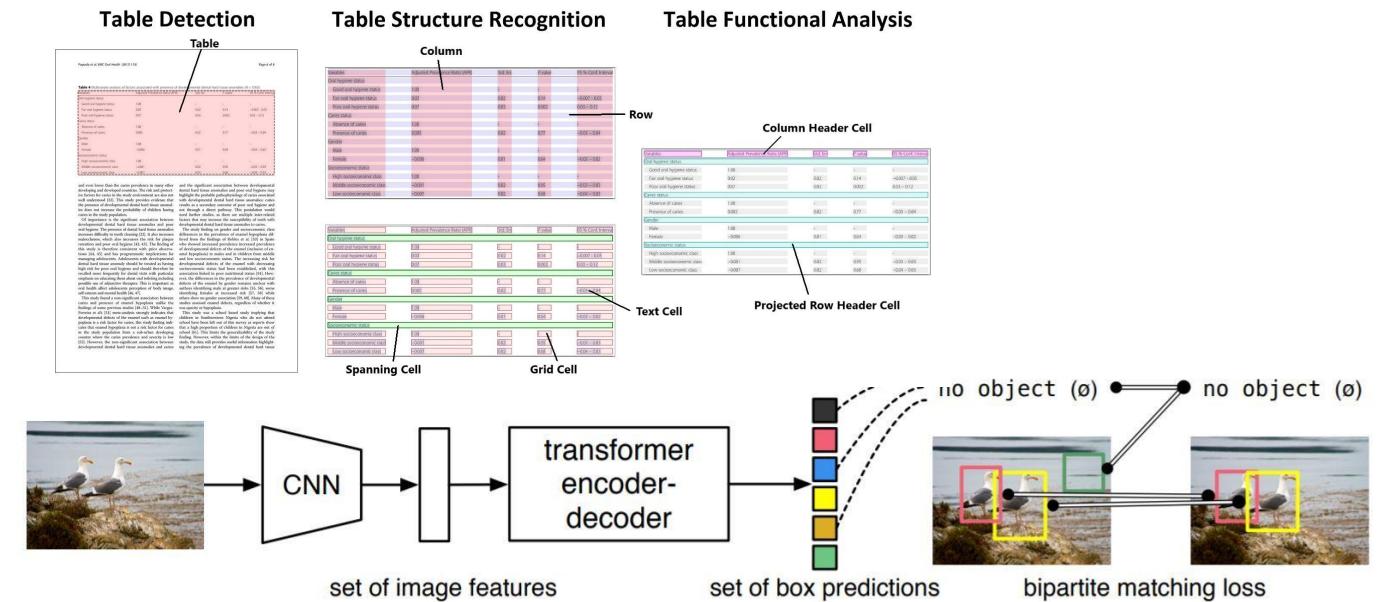
UnifiedSKG: <https://arxiv.org/abs/2201.05966>  
TaPEX: <https://arxiv.org/abs/2107.07653>

# Vision Encoder Decoder

- Multimodality
- Tableformer [CVPR22]: Table detection & OCR
  - Bounding box detection, structure generation



- TATR [CVPR22]
  - Table detection
  - Based on object detection transformer (DETR)

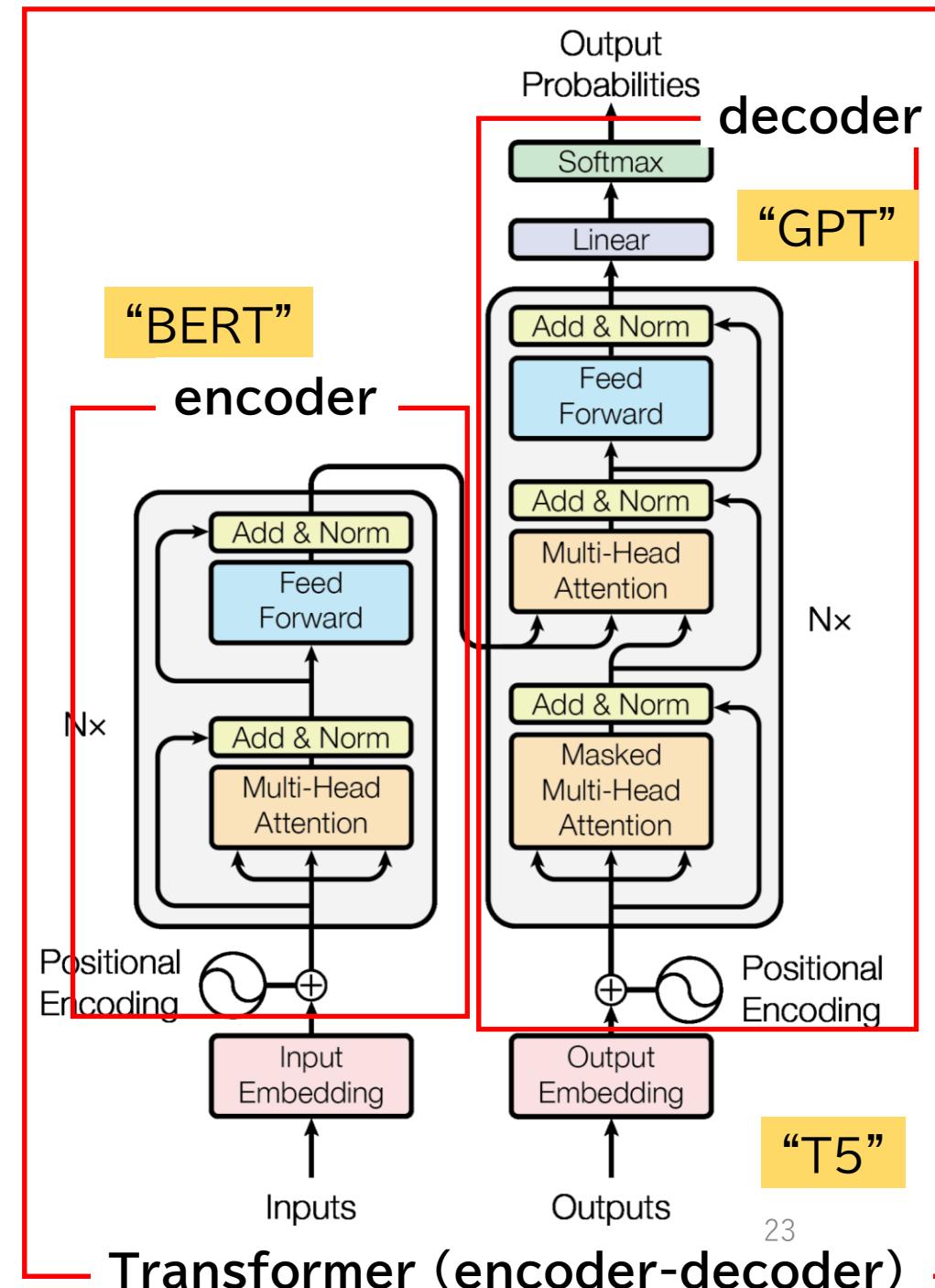


Tableformer: <https://arxiv.org/abs/2203.01017>  
 TATR: <https://github.com/microsoft/table-transformer>

TATR(DETR)

# LLM (decoders for table)

- Autoregressive Generation
  - Good for generative task
  - Easy for self-supervised training
- Simple architecture
  - Efficient on training large-scale data
  - Easy to scale the model



# Motivation of LLM for table task

- Good points
  - Large knowledge
  - In context learning/ zero shot ability
  - Generation task, interactive
- Limitation
  - Slow
  - Cost
  - Not good performance for
    - Mathematical task
    - Large table

# Use of LLM in table tasks

Prompt

Fine-tuning

RAG

LLM agents

VLM

# On the Use of Large Language Models for Table Tasks

- Prompting

# Prompt Engineering

## Zero-shot

- Instruct the LLM to follow prompts.

## Few-shot

- Demonstrate with one or more examples.

## Chain-of-thought

- Elicit complex reasoning through a series of intermediate reasoning steps.
- Variants: **tree-of-thought**, **graph-of-thought**, etc.



## Bad Prompt

I'm hungry.



Certainly, here's a whole bunch of useless text because I don't know what you want because I can't read your mind because all I do is generate text from prompts. If you are unsatisfied with my response please provide a more useful prompt.



## Good Prompt

I'm an adult man weighing 200 pounds with a daily caloric requirement of 2500. I need to hit my protein goal for the day which is 120 grams. I'm hungry for something hot and a little salty.



Scramble and cook 4 eggs.

source: <https://medium.com/the-augmented-engineer/new-to-ai-prompt-writing-learn-how-to-get-summaries-that-cover-all-the-key-points-you-care-about-46615fab135b>

# Zero-Shot Prompting

- Instruct the model to follow specific prompts.

I went to the market and bought 10 apples. I gave two to my neighbor and two to the repairman. Then I bought 5 apples and ate one. How many apples are left?

Let's think about this step by step.



First, I bought 10 apples.



1. I gave two apples to my neighbor;  $10 - 2 = 8$  apples left.
2. I gave two apples to the repairman;  $8 - 2 = 6$  apples left.
3. Then I bought 5 apples.  $6 + 5 = 11$  apples left.
4. One apple eaten;  $11 - 1 = 10$  apples left.

As a result, there are 10 apples left.

# Few-Shot Prompting

- Demonstrate task-solving with examples provided in prompts.

Fantastic! // Positive

Terrible! // Negative

The movie was amazing! // Positive

What a horrible show! //



Negative

source: <https://qiita.com/Jazuma/items/9de7c56fd9fa8b9c1648>

# Chain-of-Thought

- Elicit complex reasoning by providing inference processes.

The odd numbers in this group add up to an even number. : 4, 8, 9, 15, 12, 2, 1.

A: Adding all odd numbers gives  $9+15+1 = 25$ . The answer is False.



The odd numbers in this group add up to an even number. : 15, 32, 5, 13, 82, 7, 1.

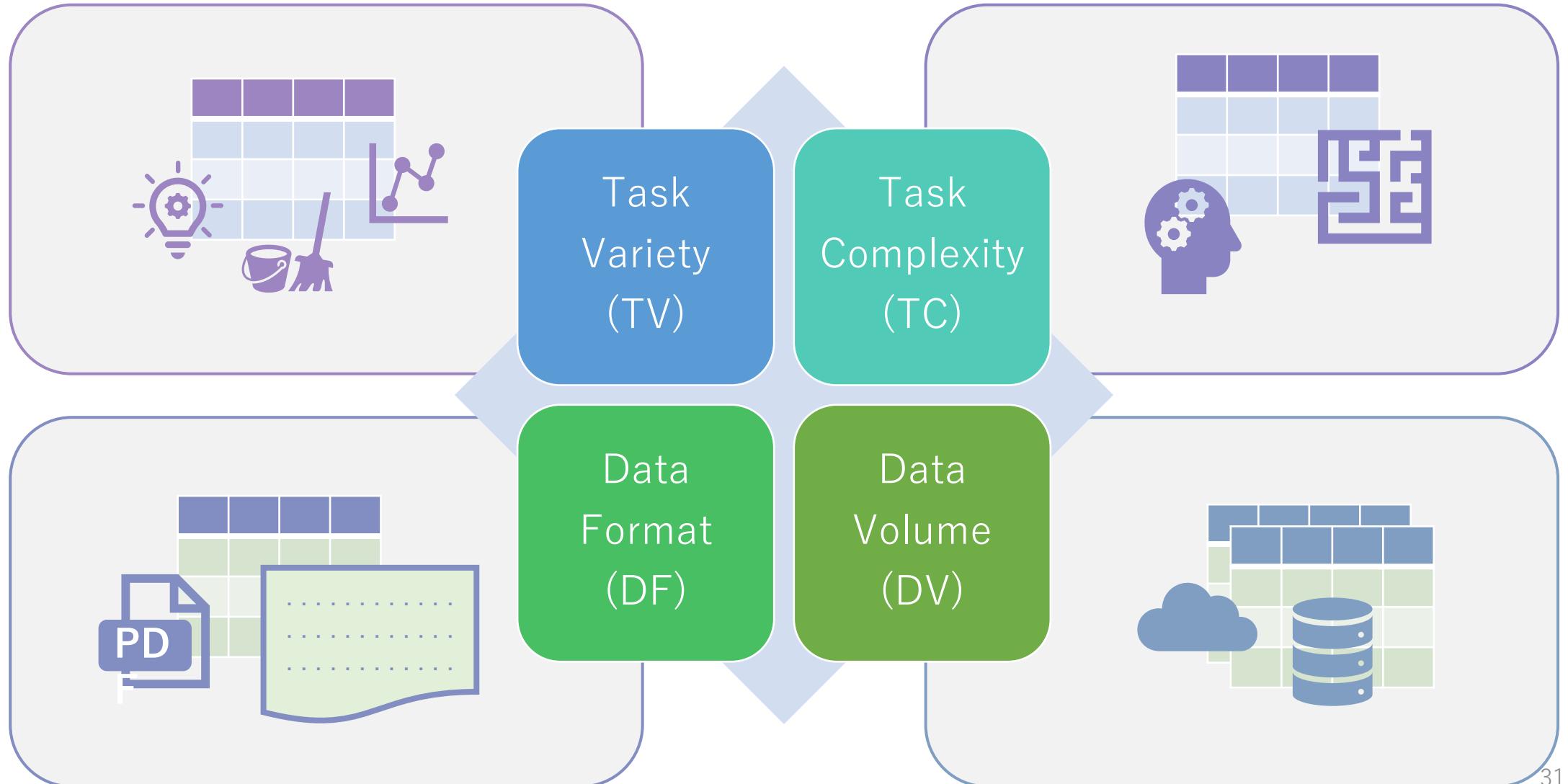
A:



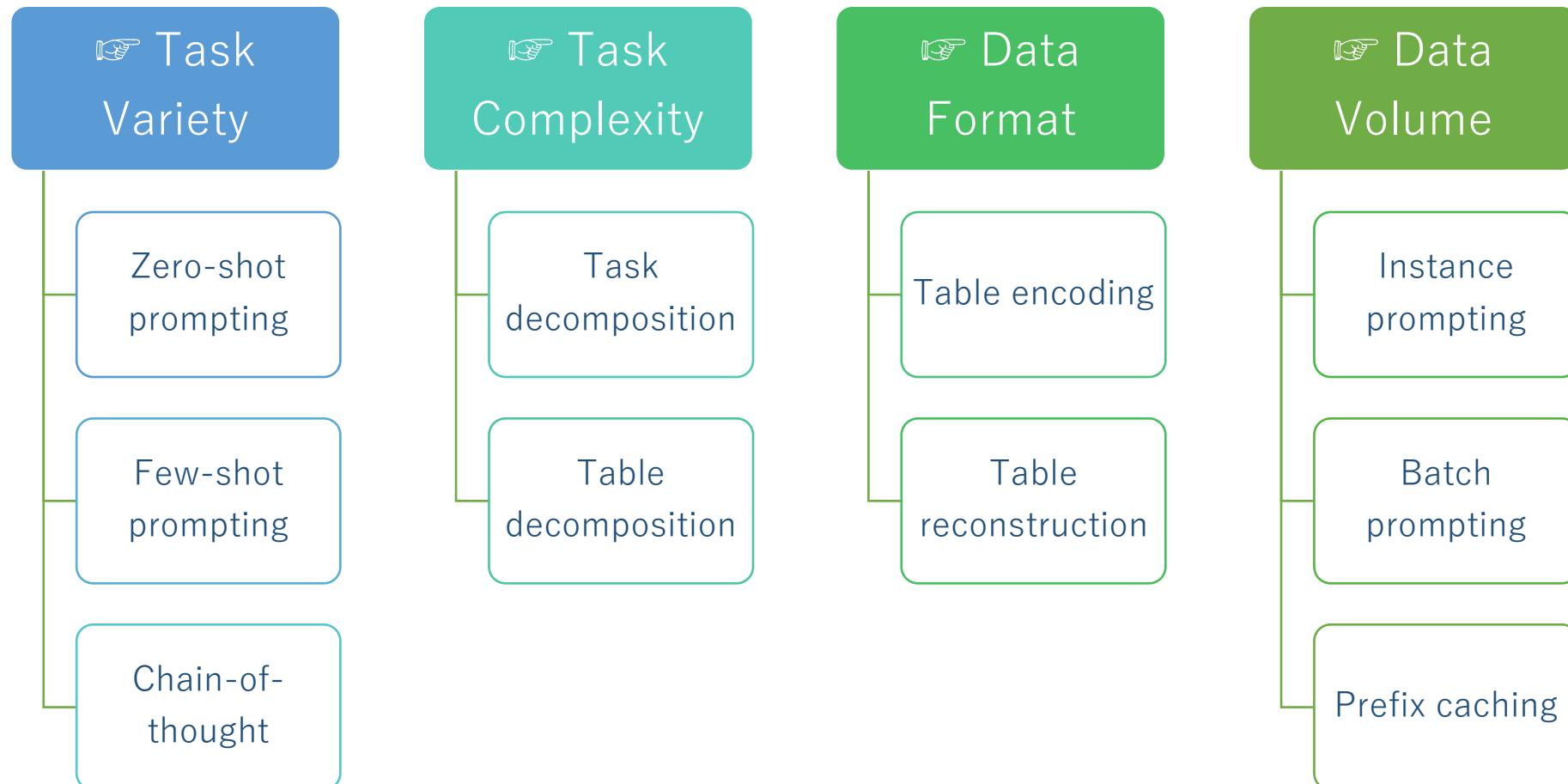
Adding all the odd numbers gives  $15 + 5 + 13 + 7 + 1 = 41$ . The answer is False.

source: <https://qiita.com/Jazuma/items/9de7c56fd9fa8b9c1648>

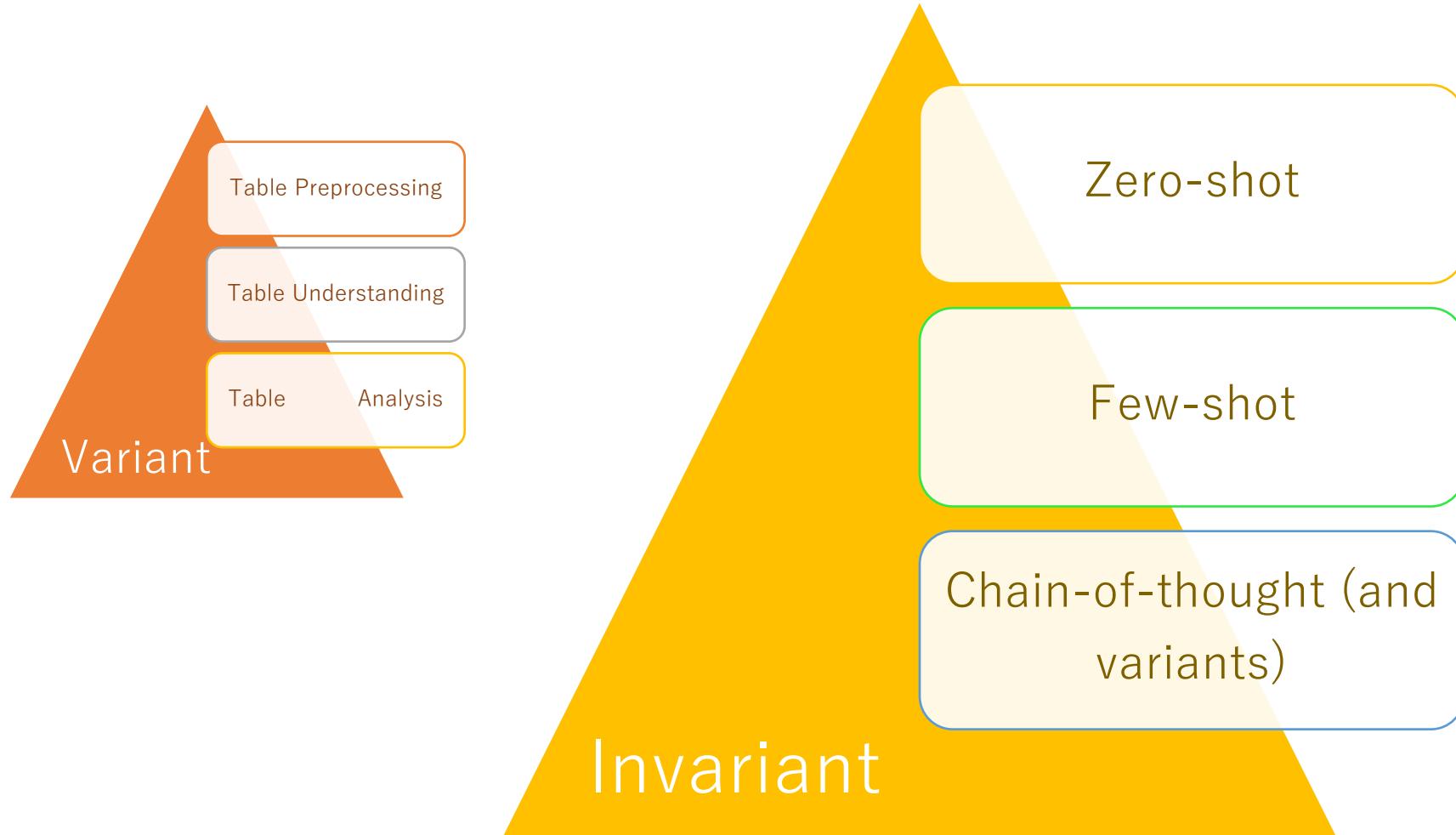
# Issues of Prompting for Table Tasks



# Prompting Techniques for Table Tasks



# 👉 Task Variety



 Task Variety: Zero-Shot Prompting

| name     | city | addr                | phone        | Type  |
|----------|------|---------------------|--------------|-------|
| Langer's | ?    | 704 S. Alvarado St. | 213-483-8050 | delis |

| Task Type        | Data Imputation  |
|------------------|--|
| Task Description | You are a database engineer.<br><br>You are requested to <b>infer the value of the "city" attribute</b> based on the values of other attributes. |
| Data Instance    | [name: "langer's", addr: "704 s. alvarado st.", phone: "213-483-8050", type: "delis"]  |

*prompt*

Answer

The city is "Los Angeles".



# Task Variety: Few-Shot Prompting

| Task Type         | Data Imputation  |
|-------------------|--|
| Task Description  | <p>You are a database engineer.</p> <p>You are requested to infer the value of the "city" attribute based on the values of other attributes.</p>   |
| Data Instance     | [name: "langer's", addr: "704 s. alvarado st.", phone: "213-483-8050", type: "delis"]  |
| Few-shot Examples | <p><b>Some examples are given below.</b></p> <p>...</p>  |
|                   | <p>User:</p> <p>Question 1: Record is [name: "carey's corner", addr: "1215 powers ferry rd.", phone: "770-933-0909", type: "hamburgers"]. What is the city?</p> <p>Assistant:</p> <p>Answer 1: Marietta</p> <p>...</p> |
|                   |  |
| Answer            | Los Angeles  |

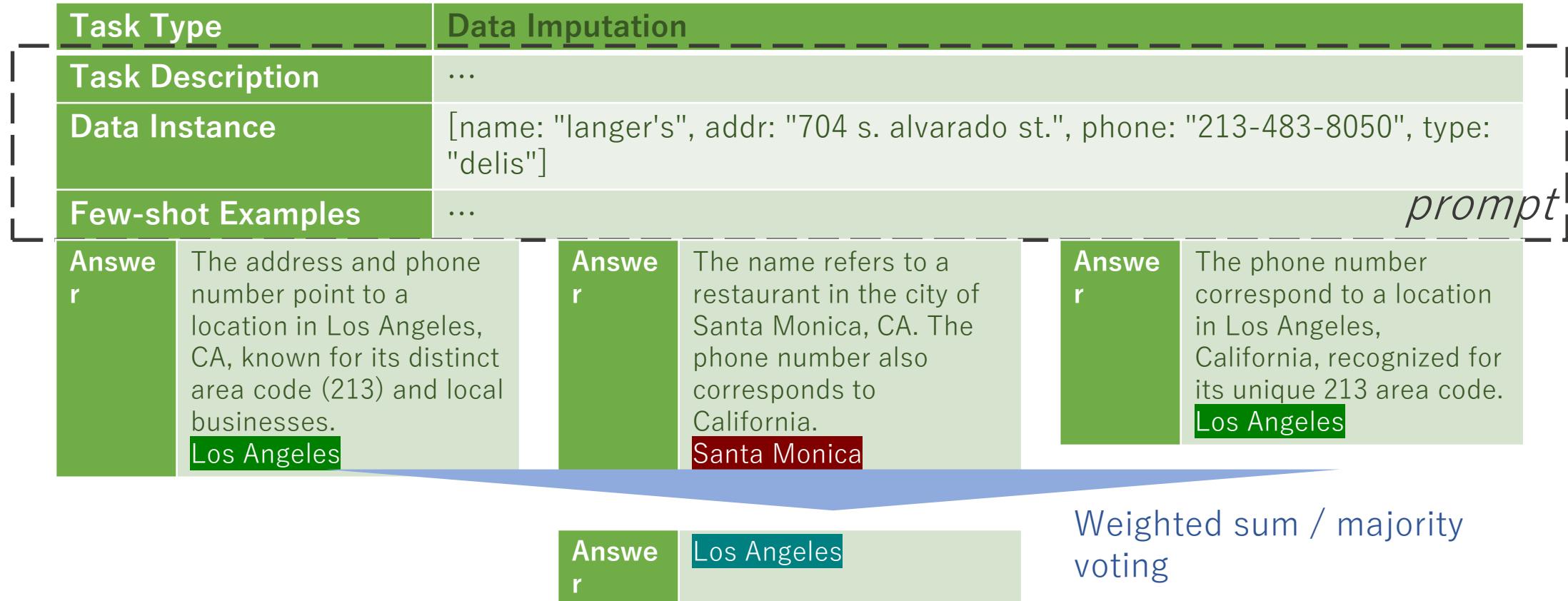


# Task Variety: Chain-of-Thought

| Task Type         | Data Imputation  |
|-------------------|--|
| Task Description  | <p>...</p> <p>MUST answer each question in two lines. In the first line, you <b>give the reason</b> for the inference. In the second line, you <b>ONLY give the value</b> of the "city" attribute.</p>   |
| Data Instance     | [name: "langer's", addr: "704 s. alvarado st.", phone: "213-483-8050", type: "delis"]  |
| Few-shot Examples | <p>Some examples are given below.</p> <p>...</p> <p>User:</p> <p>Question 1: Record is [name: "carey's corner", addr: "1215 powers ferry rd.", phone: "770-933-0909", type: "hamburgers"]. What is the city?</p> <p>Assistant:</p> <p>Answer 1: The phone area code 770 corresponds to the Atlanta metropolitan area in Georgia, and "1215 Powers Ferry Rd." is an address located in Marietta, Georgia; therefore, the city is Marietta.</p> <p>Marietta</p> <p><i>prompt</i></p> |
| Answer            | <p>The address and phone number point to a location in Los Angeles, CA, known for its distinct area code (213) and local businesses.</p> <p>Los Angeles</p>  |

# 👉 Task Variety: Chain-of-Thought (+Self-Consistency Decoding)

- Generate multiple reasoning answers and aggregate them.



Chen. Large language models are few (1)-shot table reasoners. EACL 2023.

Chen et al. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. TMLR 2023.

TA → Table  
QA

# 👉 Task Complexity



# 👉 Task Complexity: Task Decomposition



| Task Type        | Column Type Annotation  |
|------------------|---|
| Task Description | Classify the columns of a given table with only one of the following classes that are separated with comma:<br><u>description of event</u> , <u>description of restaurant</u> , <u>postal code</u> , <u>region of address</u> ... |
| Data Instance    | Column 1    Column 2    Column 3    Column 4<br>Friends Pizza    2525    Cash Visa MasterCard    7:30AM   |
| Answer           | name, number, payment, time   |

**Korini and Bizer. Column type annotation using ChatGPT. TaDA 2023.**

TU → Table Interpretation

Zhao et al. Large language models are complex table parsers. EMNLP 2023.

TA → Table QA

Wang et al. Chain-of-table: Evolving tables in the reasoning chain for table understanding. ICLR 2024.

TP → Table Transformation

TA → Table QA

Dong et al. OpenTE: Open-structure table extraction from text. ICCASP 2024.

TU → Table Interpretation



# Task Complexity: Task Decomposition

Sub-Task 1

After decomposition

Sub-Task 2

| Sub-Task Type    | Table Classification   |
|------------------|--|
| Task Description | Your task is to <b>classify if a table</b> describes <u>Restaurants</u> , <u>Events</u> , <u>Music Recordings</u> , or <u>Hotels</u> . |
| Data Instance    | Column 1    Column 2    Column 3   <br>Column 4 ¶<br>Friends Pizza    2525    Cash Visa<br>MasterCard    7:30AM ¶                      |
| Answer           | Restaurant   |

| Sub-Task Type    | Column Classification   |
|------------------|---|
| Task Description | Your task is to <b>classify the columns</b> of a given table with only one of the following classes that are separated with comma: <u>name of restaurant</u> , <u>description of restaurant</u> ... |
| Data Instance    | Column 1    Column 2    Column 3   <br>Column 4 ¶<br>Friends Pizza    2525    Cash Visa<br>MasterCard    7:30AM ¶   |
| Answer           | name of restaurant, postal code, payment accepted, time   |

**Korini and Bizer. Column type annotation using ChatGPT. TaDA 2023.**

Zhao et al. Large language models are complex table parsers. EMNLP 2023.

Wang et al. Chain-of-table: Evolving tables in the reasoning chain for table understanding. ICLR 2024.

Dong et al. OpenTE: Open-structure table extraction from text. ICCASP 2024.

TU → Table Interpretation

TA → Table QA

TP → Table Transformation

TA → Table QA

TU → Table Interpretation

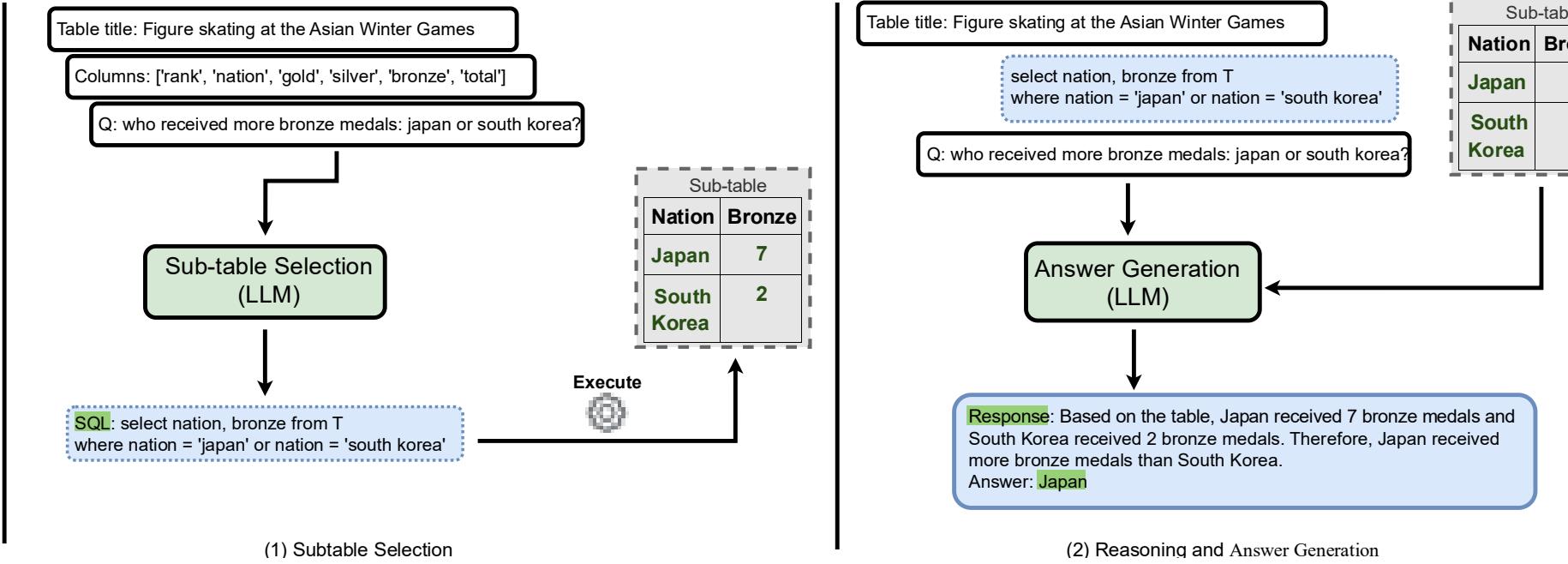


# Task Complexity: Table Decomposition

Table: Figure skating at the Asian Winter Games

| Rank  | Nation      | Gold | Silver | Bronze | Total |
|-------|-------------|------|--------|--------|-------|
| 1     | China       | 13   | 9      | 13     | 35    |
| 2     | Japan       | 7    | 10     | 7      | 24    |
| 3     | Uzbekistan  | 1    | 2      | 3      | 6     |
| 4     | Kazakhstan  | 2    | 2      | 0      | 4     |
| 5     | North Korea | 1    | 0      | 1      | 2     |
| 6     | South Korea | 0    | 0      | 2      | 2     |
| Total |             | 24   | 23     | 26     | 73    |

Q: who received more bronze medals: japan or south korea?  
A: Japan



**Nahid and Rafiei. TabSQLify: Enhancing reasoning capabilities of LLMs through table decomposition. NAACL 2024.**

Patnaik et al. Cabinet: Content relevance based noise reduction for table question answering. ICLR 2024.

Jiang et al. StructGPT: A general framework for large language model to reason over structured data. EMNLP 2023.

TA → Table QA

TA → Text-to-SQL

TA → Table QA

TA → Table QA

TA → Text-to-SQL

# 👉 Task Complexity: Table Decomposition (Progressive Prompting)

| Question        | Report the number of wins in Grand Slam tournaments. |       |             |                 |     |       |             | Text transformation |     |       |             |        |
|-----------------|--|-------|-------------|-----------------|-----|-------|-------------|---------------------|-----|-------|-------------|--------|
| tournament      | ...  | att n | career_w /l | tournament      | ... | att n | career_w /l | tournament          | ... | att n | career_w /l | n_wi n |
| Australian Open | ...  | 18    | 22-18       | Australian Open | ... | 18    | 22-18       | Australian Open     | ... | 18    | 22-18       | 22     |
| Roland Garros   | ...  | 14    | 11-14       | Roland Garros   | ... | 14    | 11-14       | Roland Garros       | ... | 14    | 11-14       | 11     |
| Wimbledon       | ...  | 18    | 13-18       | Wimbledon       | ... | 18    | 13-18       | Wimbledon           | ... | 18    | 13-18       | 13     |
| US Open         | ...  | 13    | 16-13       | US Open         | ... | 13    | 16-13       | US Open             | ... | 13    | 16-13       | 16     |
| Indian Wells    | ...  | 15    | 20-15       | Indian Wells    | ... | 15    | 20-15       | Indian Wells        | ... | 15    | 20-15       |        |
| ...             | ...  | ...   | ...         | ...             | ... | ...   | ...         | ...                 | ... | ...   | ...         |        |

Focus on **column selection**.

Set the groundwork for understanding how to fetch specific data from a database.

Incorporate both **column and row selection**.

Extract particular columns and filtering rows based on specified criteria, enhancing precision in data gathering.

Apply **additional operations** (e.g., aggregation functions and text operations).

Aggregation functions empower data summarization.  
Text operations facilitate the manipulation and transformation of string data.

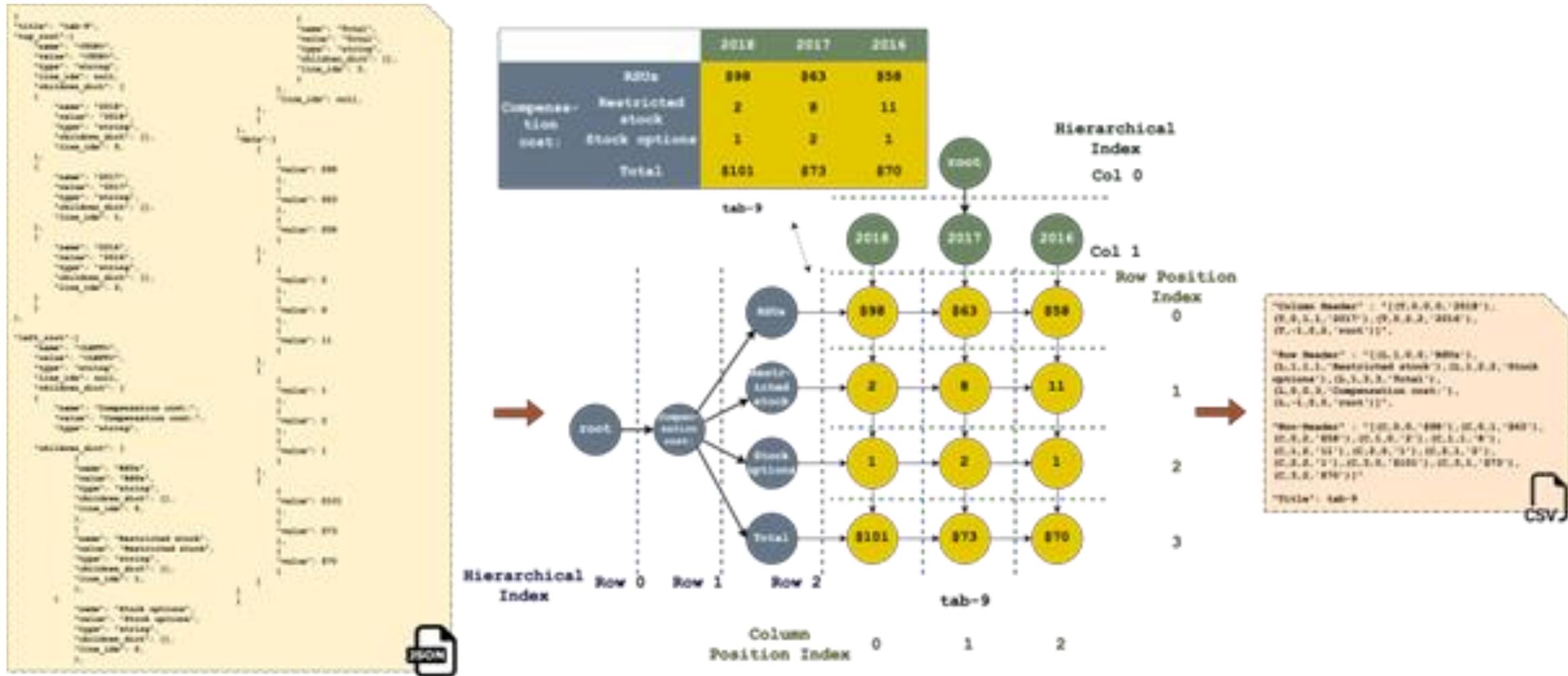


# Data Format: Table Encoding

| text<br>(serialized) | spreadsheet | markup | key-value | program  | image   | embedded |
|----------------------|-------------|--------|-----------|----------|---------|----------|
| <br>TXT              | <br>CSV     | <br>MD | <br>HTML  | <br>JSON | <br>SQL | <br>JPG  |

easy hard

# Data Format: Table Reconstruction



# 👉 Data Volume: Instance Prompting

| name       | city | addr                | phone        | Type    |
|------------|------|---------------------|--------------|---------|
| Langer's   | ?    | 704 S. Alvarado St. | 213-483-8050 | delis   |
| Valetino   | ?    | 3115 Pico Blvd.     | 310-829-4313 | Italian |
| Cafe Bizou | ?    | 14016 Ventura Blvd. | 818/788-3536 | French  |

| Task Type        | Data Imputation   | Task Type        | Data Imputation   | Task Type        | Data Imputation   |
|------------------|---|------------------|---|------------------|---|
| Task Description | You are a database engineer.<br><br>You are requested to infer the value of the "city" attribute based on the values of other attributes. | Task Description | You are a database engineer.<br><br>You are requested to infer the value of the "city" attribute based on the values of other attributes. | Task Description | You are a database engineer.<br><br>You are requested to infer the value of the "city" attribute based on the values of other attributes. |
| Data Instance    | [name: "langer's", addr: "704 s. alvarado st.", phone: "213-483-8050", type: "delis"]<br><i>prompt</i>                                    | Data Instance    | [name: "valentino", addr: "3115 pico blvd.", phone: "310-829-4313", type: "italian"]<br><i>prompt</i>                                     | Data Instance    | [name: "cafe bizou", addr: "14016 ventura blvd.", phone: "818/788-3536", type: "french"]<br><i>prompt</i>                                 |



# Data Volume: Batch Prompting

| name       | city | addr                | phone        | Type    |
|------------|------|---------------------|--------------|---------|
| Langer's   | ?    | 704 S. Alvarado St. | 213-483-8050 | delis   |
| Valetino   | ?    | 3115 Pico Blvd.     | 310-829-4313 | Italian |
| Cafe Bizou | ?    | 14016 Ventura Blvd. | 818/788-3536 | French  |

| Task Type        | Data Imputation   |
|------------------|---|
| Task Description | You are a database engineer.<br><br>You are requested to infer the value of the "city" attribute based on the values of other attributes.   |
| Data Instance    | [name: "langer's", addr: "704 s. alvarado st.", phone: "213-483-8050", type: "delis"]<br>[name: "valentino", addr: "3115 pico blvd.", phone: "310-829-4313", type: "italian"]<br>[name: "cafe bizou", addr: "14016 ventura blvd.", phone: "818/788-3536", type: "french"] |

*prompt*

**Standard Prompting**

# K-shot in-context exemplars

Q: {question}  
A: {answer}

Q: {question}  
A: {answer}

...

# One sample to inference

Q: Ali had \$21. Leila gave him half of her \$ 100. How much does Ali have now?

# Response

A: Leila gave  $100/2=50$  to Ali. Ali now has  $\$21+\$50 = \$71$ . The answer is 71.

**Batch Prompting**

# K-shot in-context exemplars in K/b batches

Q[1]: {question}  
Q[2]: {question}  
A[1]: {answer}  
A[2]: {answer}

} b(=2) samples in one batch

# b samples in a batch to inference

Q[1]: Ali had \$21. Leila gave him half of her \$100. How much does Ali have now?  
Q[2]: A robe takes 2 bolts of blue fiber and half that white fiber. How many bolts?

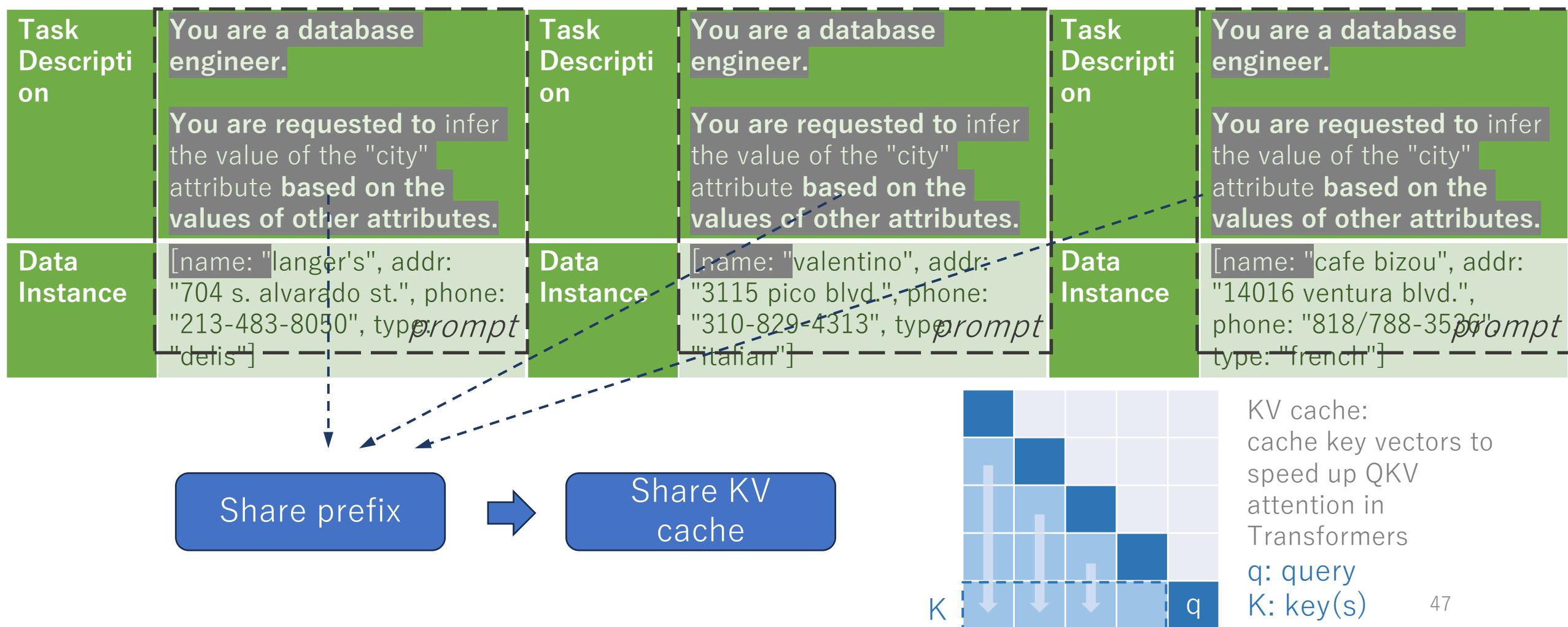
# Responses to a batch

A[1]: Leila gave  $100/2=50$  to Ali. Ali now has  $\$21+\$50 = \$71$ . The answer is 71.  
A[2]: It takes  $2/2=1$  bolt of white fiber. The total amount is  $2+1=3$ . The answer is 3.

source: Cheng et al. Batch prompting: Efficient inference with large language model APIs. EMNLP 2023.

# 👉 Data Volume: Prefix Caching (+Instance Prompting)

- Use Automatic Prefix Caching (APC) in the vLLM library.



# Summary

- 
- Prompt engineering is effective in improving the performance of table tasks.
  - Zero-shot, few-shot, and chain-of-thought are useful across a variety of tasks.
  - Task/table decomposition is recommended when dealing with complex tasks.
  - Tables can be encoded in various formats. Reconstruction may help reduce difficulty.
  - To speed up local deployment, users are suggested to use APC in the vLLM library.

# Q&A

\Orchestrating a brighter world

**NEC**

- On the Use of Large Language Models for Table Tasks
  - Prompting



**大阪大学**  
OSAKA UNIVERSITY

Orchestrating a brighter world

**NEC**



# On the Use of Large Language Models for Table Tasks

- Fine-Tuning

# A Brief Tour of LLM

# Anatomy of an Open-Source LLM

- Configuration Files

Model architecture hyperparameters  
Default settings for text generation

- Model Weight Files

Model parameters  
Index mapping model components to files

- Tokenizer Files:

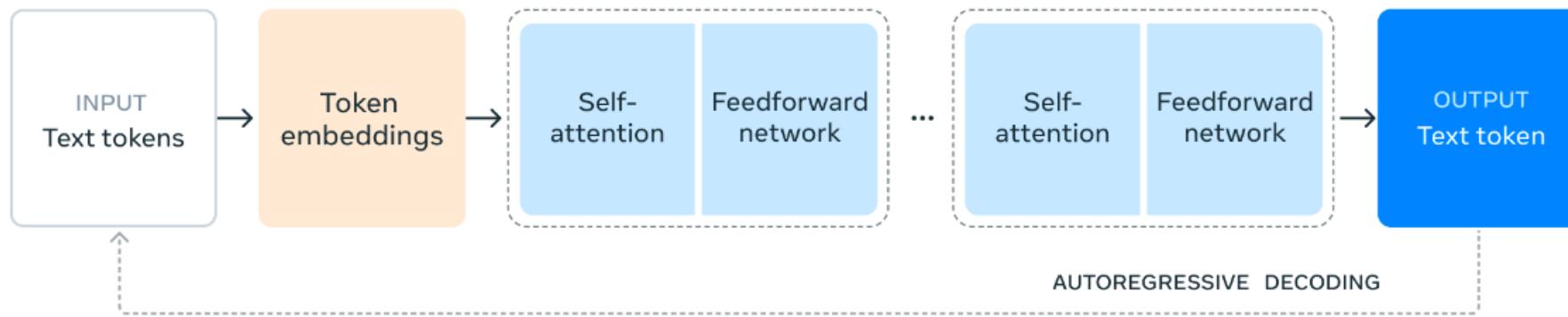
Special token maps (e.g. start of text indicator)  
Vocabulary File  
Tokenizer settings

|                                  |             |
|----------------------------------|-------------|
| LICENSE                          | 7.63 kB     |
| README.md                        | 40.9 kB     |
| USE_POLICY.md                    | 4.69 kB     |
| config.json                      | 826 Bytes   |
| generation_config.json           | 185 Bytes   |
| model-00001-of-00004.safetensors | 4.98 GB LFS |
| model-00002-of-00004.safetensors | 5 GB LFS    |
| model-00003-of-00004.safetensors | 4.92 GB LFS |
| model-00004-of-00004.safetensors | 1.17 GB LFS |
| model.safetensors.index.json     | 24 kB       |
| special_tokens_map.json          | 73 Bytes    |
| tokenizer.json                   | 9.09 MB     |
| tokenizer_config.json            | 50.5 kB     |

1. Composition of Llama 3.1 8B Model

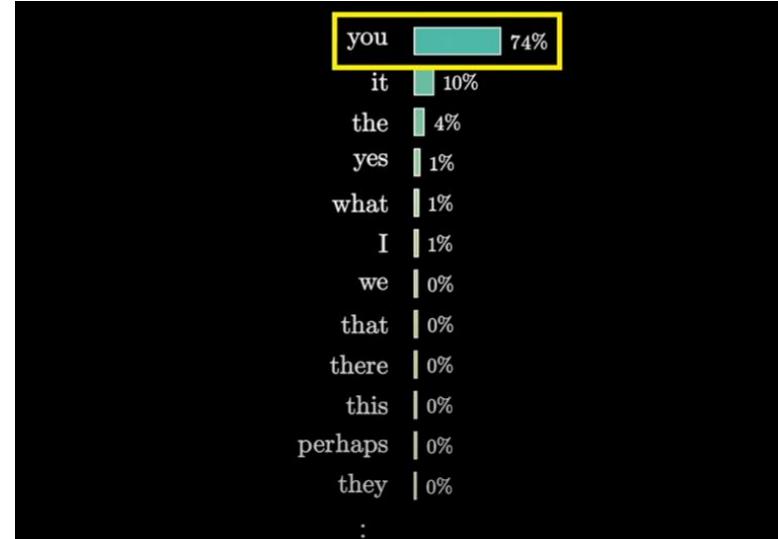
1: Meta: <https://huggingface.co/meta-llama/Llama-3.1-8B/tree/main>

# How LLMs Work



1. Typical Generation Process of Decoder-Only Transformer Models (Llama 3.1)

If you could see the underlying probability distributions a large language model uses when generating text, then **you**



2. High-Level Visualization of Autoregressive Decoding

1: Meta <https://ai.meta.com/blog/meta-llama-3-1/>

2: 3Blue1Brown

<https://www.3blue1brown.com/topics/neural-networks>

# Dive into Forward Pass



LLMs compute the logits for the next token based on input text<newline>

[7454, 25153, 23864, 290, 148063, 395, 290, 2613, 6602, 4122, 402, 3422, 2201, 198]

## 1. GPT-4o Tokenization Demonstration

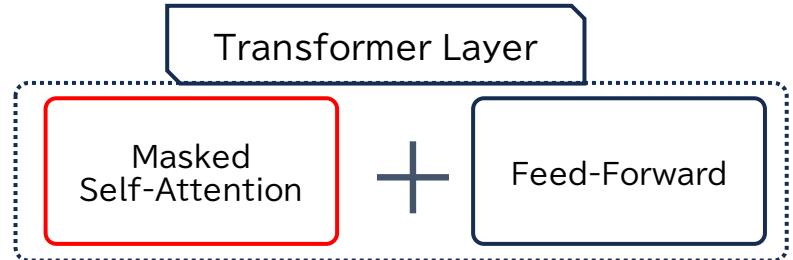
$d_{model}$   
 $N_{vocab}$

```
"hidden_size": 4096,  
"vocab_size": 128256  
"weight_map": {  
    "lm_head.weight": "model-00004-of-00004.safetensors",  
    "model.embed_tokens.weight": "model-00001-of-00004.safetensors",
```

## Rotary Position Embeddings (RoPE)

```
"rope_scaling": {  
    "factor": 8.0,  
    "low_freq_factor": 1.0,  
    "high_freq_factor": 4.0,  
    "original_max_position_embeddings": 8192,  
    "rope_type": "llama3"  
},  
"rope_theta": 500000.0,
```

# Dive into Forward Pass



- Self-Attention

$$Q = XW_Q, \quad K = XW_K, \quad V = XW_V$$

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

- Masked Self-Attention

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} + M \right) V$$

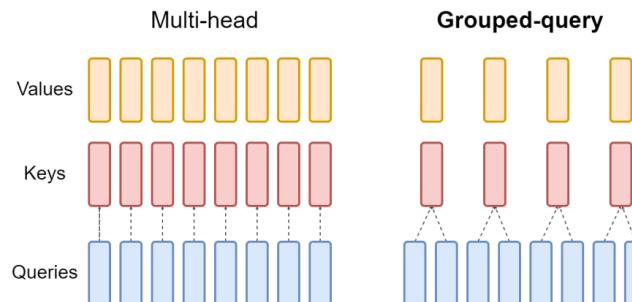
$$M_{i,j} = \begin{cases} 0, & j \leq i \\ -\infty, & j > i \end{cases}$$

- Multi-Head Attention

$$\text{Output} = \text{Concat}(\text{head}_i) W_O$$

- Each attention head computes a portion
- The outputs are concatenated

- Grouped-Query Attention



1. Overview of grouped-query method

$W_K$  "model.layers.0.self\_attn.k\_proj.weight"

$W_O$  "model.layers.0.self\_attn.o\_proj.weight"

$W_Q$  "model.layers.0.self\_attn.q\_proj.weight"

$W_V$  "model.layers.0.self\_attn.v\_proj.weight"

"num\_attention\_heads": 32,

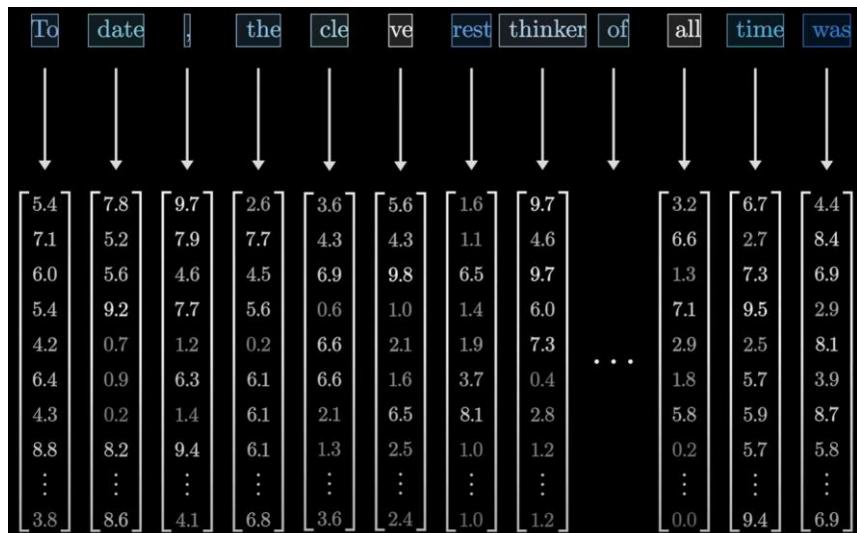
"num\_hidden\_layers": 32,

"num\_key\_value\_heads": 8,

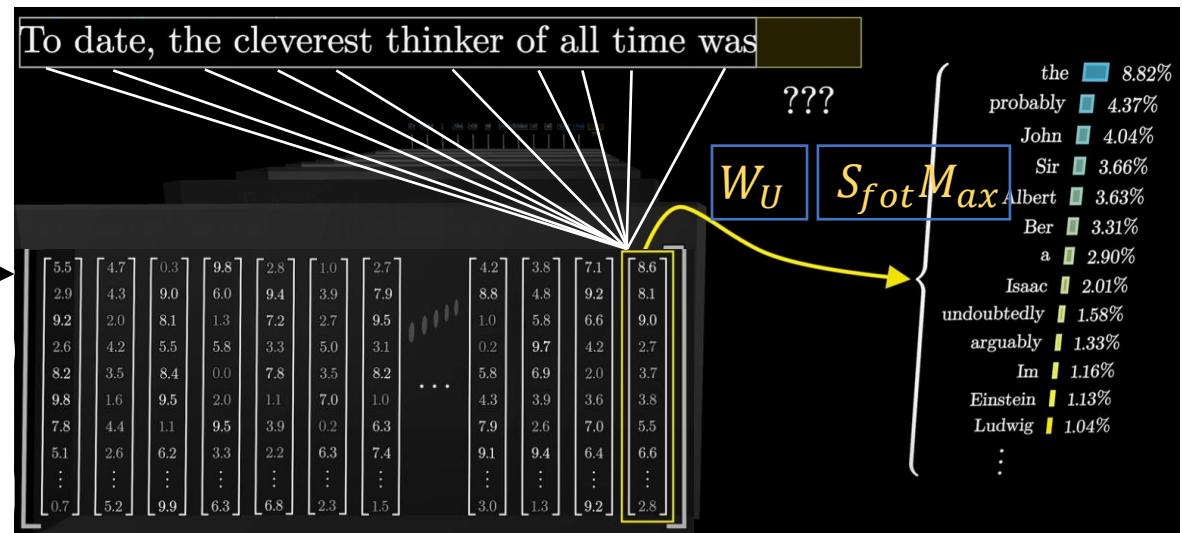
# Dive into Forward Pass

Output Layer/ Decoding

- Unembedding Matrix:  $W_U$



$W_U$  "lm\_head.weight": "model-00004-of-00004.safetensors"



1. Comparison between Input Embeddings and Final Outputs

- Tied/Untied Input-output Embeddings

- If the unembedding matrix is the transformation of the embedding matrix or a separate matrix

"tie\_word\_embeddings": `false`,

# Dive into Forward Pass

## Output Layer/ Decoding

- Greedy Decoding

Always selecting the option with the highest probability

- Temperature

Controls the **sharpness** or **smoothness** of the distribution

$$P_i = \frac{\exp(\text{logit}_i/T)}{\sum_j \exp(\text{logit}_j/T)}$$

In practice, greedy decoding is used when T= 0.

- Top-p Sampling (Nucleus Sampling)

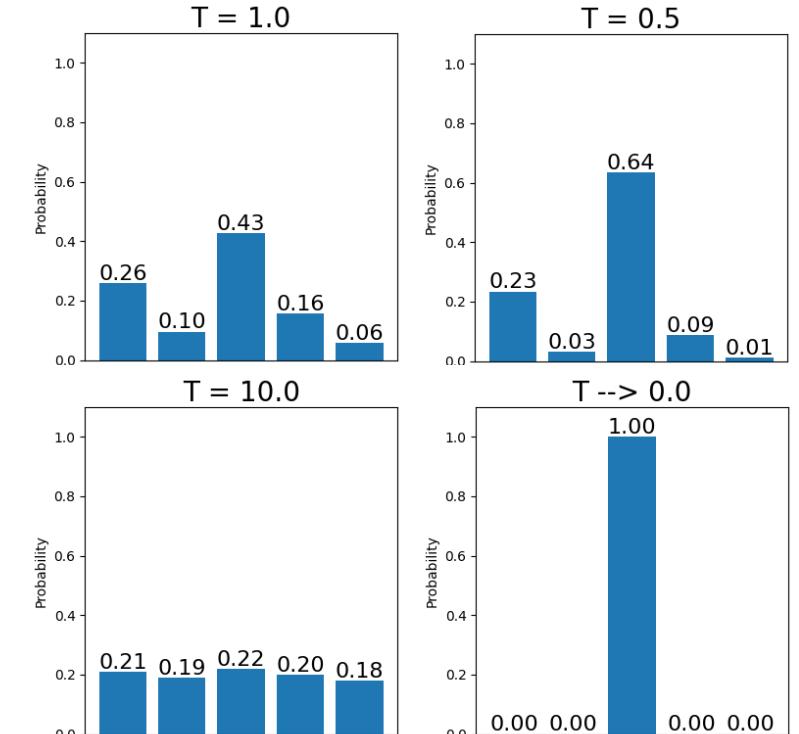
- Cumulative Probability
- Nucleus (Token Set) at Threshold p
- Re-Normalization within the Nucleus

- Top-k Sampling, Beam Search, etc.

"do\_sample": true,  
"temperature": 0.6,  
"top\_p": 0.9,

### 1. Default Decoding Settings for Llama3.1 8B

*Logits: [1.5, 0.5, 2.0, 1.0, 0.0]*



### 2. Demonstration of SoftMax with Temperature

# Recap on Llama3.1 8B

```

"attention_bias": false, Context Window "max_position_embeddings": 131072, WUp
"attention_dropout": 0.0, "mlp_bias": false, WGate
"bos_token_id": 128000, "model_type": "llama", WDown
"eos_token_id": 128001, "num_attention_heads": 32,
"hidden_act": "silu", Nheads "num_hidden_layers": 32,
dmodel "hidden_size": 4096, Natt.layers "num_key_value_heads": 8,
"initializer_range": 0.02, "pretraining_tp": 1,
dFF "intermediate_size": 14336, "rms_norm_eps": 1e-05,
dk = dmodel/Nheads

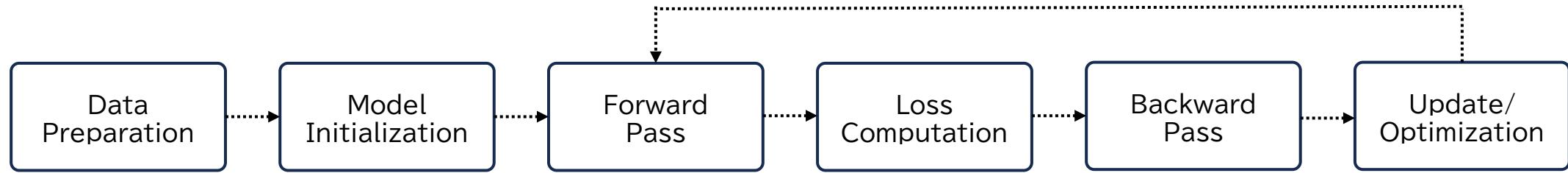
```

## 1. Llama 3.1-8B Configuration Breakdown

|                   |                      |                  |            | Per Layer   | Layers            |                           |
|-------------------|----------------------|------------------|------------|---|-------------------|---------------------------|
| Transformer Layer | Multi-head Attention | Embedding Matrix | $W_E$      | $N_{vocab} * d_{model} = 128,256 * 4,096 = 525,336,576$       | 1                 | 525,336,576               |
|                   |                      | Query Matrix     | $W_Q$      | $d_{model} * d_k * N_{heads} = 4,096 * 128 * 32 = 16,777,216$ |                   | 6,979,321,856             |
|                   |                      | Key Matrix       | $W_K$      | $d_{model} * d_k * N_{kvheads} = 4,096 * 128 * 8 = 4,194,304$ |                   |                           |
|                   |                      | Value Matrix     | $W_V$      | $d_{model} * d_k * N_{kvheads} = 4,096 * 128 * 8 = 4,194,304$ |                   |                           |
|                   | MLP                  | Output Matrix    | $W_O$      | $d_{model} * d_{model} = 4096 * 4,096 = 16,777,216$           | $N_{layers} = 32$ | 6,979,321,856             |
|                   |                      | Up Projection    | $W_{Up}$   | $d_{model} * d_{FF} = 4096 * 14,336 = 58,720,256$             |                   |                           |
|                   |                      | Gate Projection  | $W_{Gate}$ | $d_{model} * d_{FF} = 4096 * 14,336 = 58,720,256$             |                   |                           |
|                   |                      | Down Projection  | $W_{Down}$ | $d_{FF} * d_{model} = 14,336 * 4096 = 58,720,256$             |                   |                           |
|                   | Unembedding Matrix   | LayerNorm        | $W_U$      | $d_{model} * N_{vocab} = 4,096 * 128,256 = 525,336,576$       | 1                 | 525,336,576               |
|                   |                      |                  | $W_{LN}$   | $d_{model} = 4,096$   |                   | $2 * N_{layers} + 1 = 65$ |
|                   |                      |                  |            |   | Total             | 8,030,261,248             |

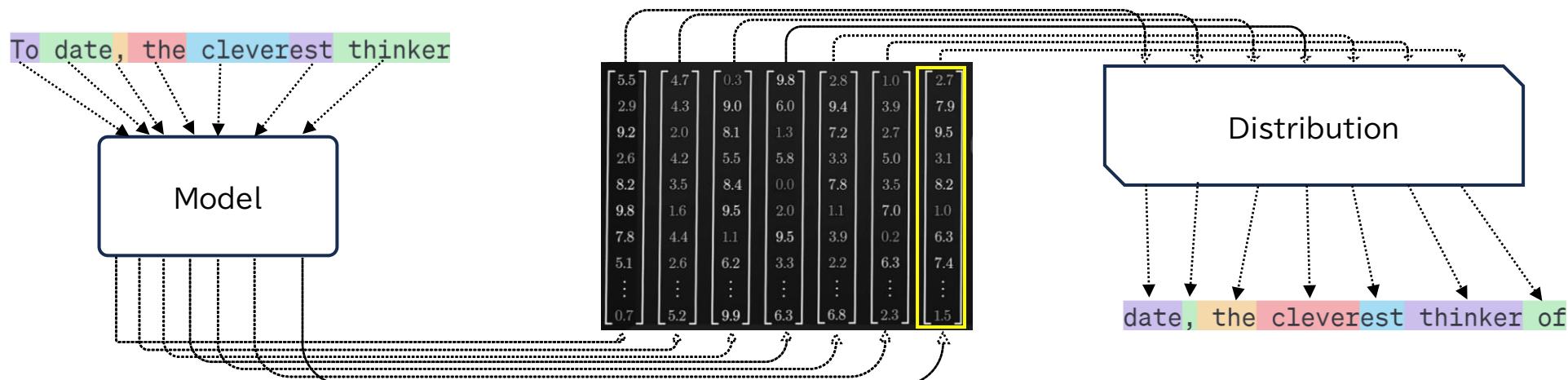
1: Meta <https://huggingface.co/meta-llama/Llama-3.1-8B/tree/main>

# Fundamentals of LLMs Training



- Teacher Forcing
  - Uses ground truth next tokens as inputs for each time step
- Target Shifting
  - Processes entire input sequence in one forward pass.

$$X = [x_1, x_2, \dots, x_L], \quad Y = [x_2, x_3, \dots, x_{L+1}]$$



# Pre-Training and Fine-Tuning

- Pre-Training

- Initialization

- Randomly initialize weights

- Data

- Large scale, diverse text corpus

- Purpose

- Learning general language representations and knowledge



The benefits of a healthy diet include



improved energy levels, better mood, and reduced risk of chronic diseases.

Llama3.1-8B

- Fine-Tuning

- Initialization

- Load pre-trained weights

- Data

- Task-specific or domain-specific

- Purpose

- Following instructions or engaging in conversations

<|user|> List the benefits of a healthy diet.  
<|endoftext|><|assistant|>



A healthy diet offers numerous benefits, including:

1. Improved energy levels

...

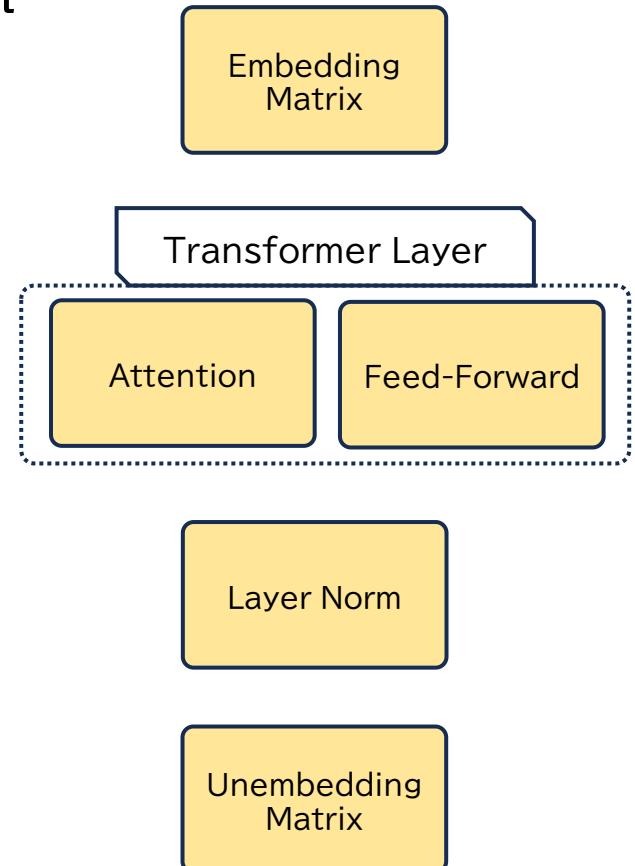


Llama-3.1-8B-Instruct

# Fine-Tuning Methods

# Full Parameter Fine-Tuning

- Adjusting **all the parameters** of a pre-trained model
- Advantages
  - Maximum Adaptability
  - Potential for Higher Performance
- Disadvantages
  - Computationally Intensive
  - Risk of Overfitting
- Use Case
  - Sufficient training data is available
  - The new task is significantly different



# Table-GPT

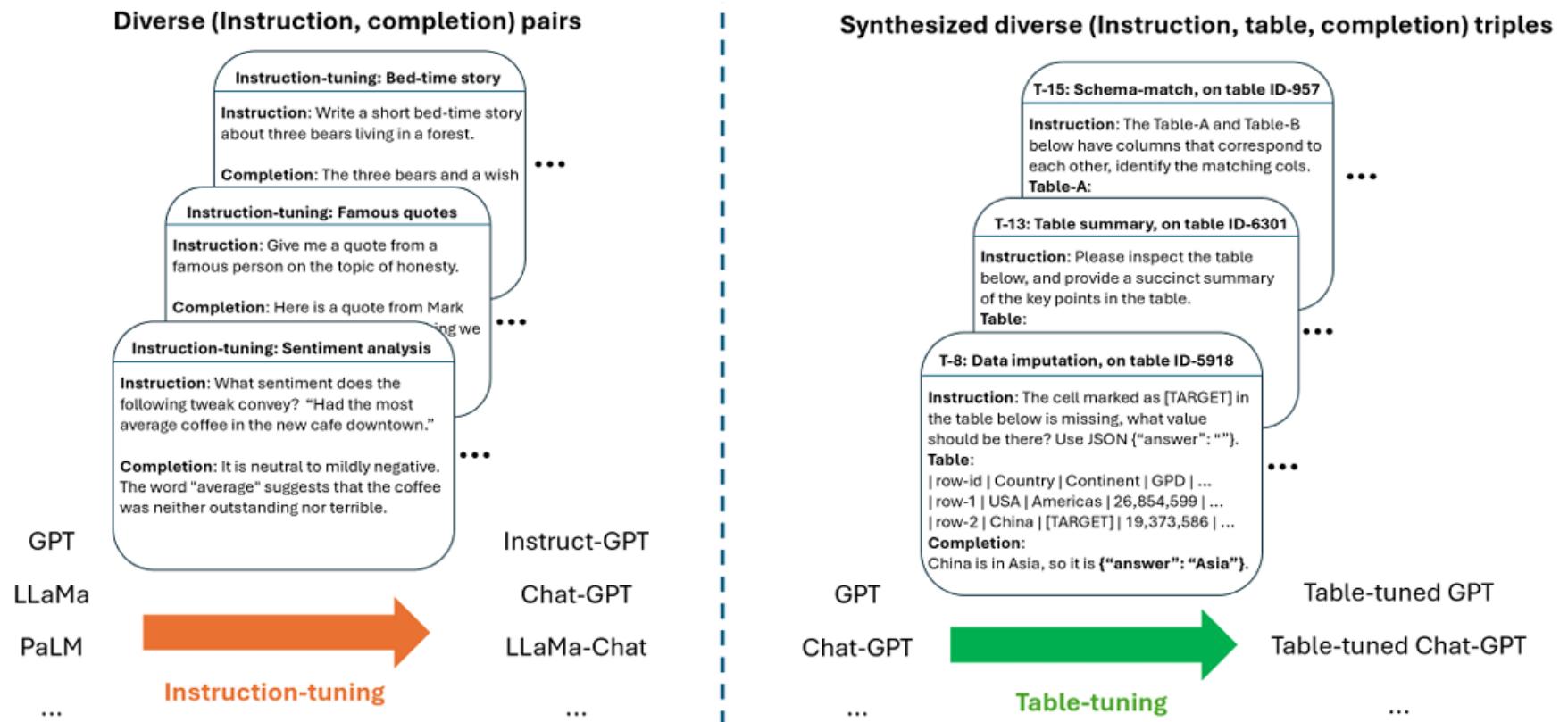
Table Preprocessing

Table Understanding

Table Analysis

GPT3.5  
ChatGPT

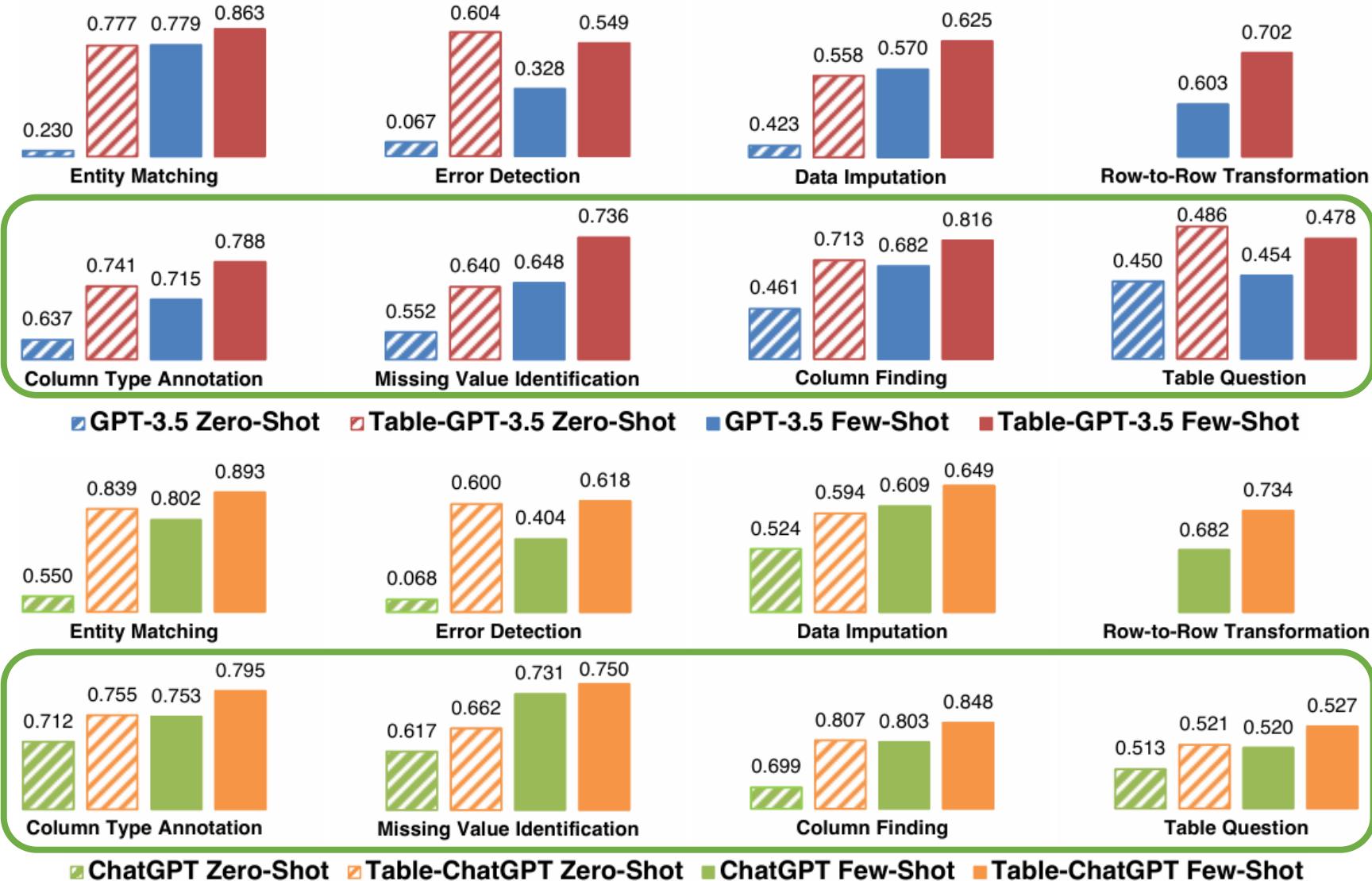
- Extensive Task Coverage
- Synthesis-then-Augment
- Comprehensive Analysis



# Table-GPT

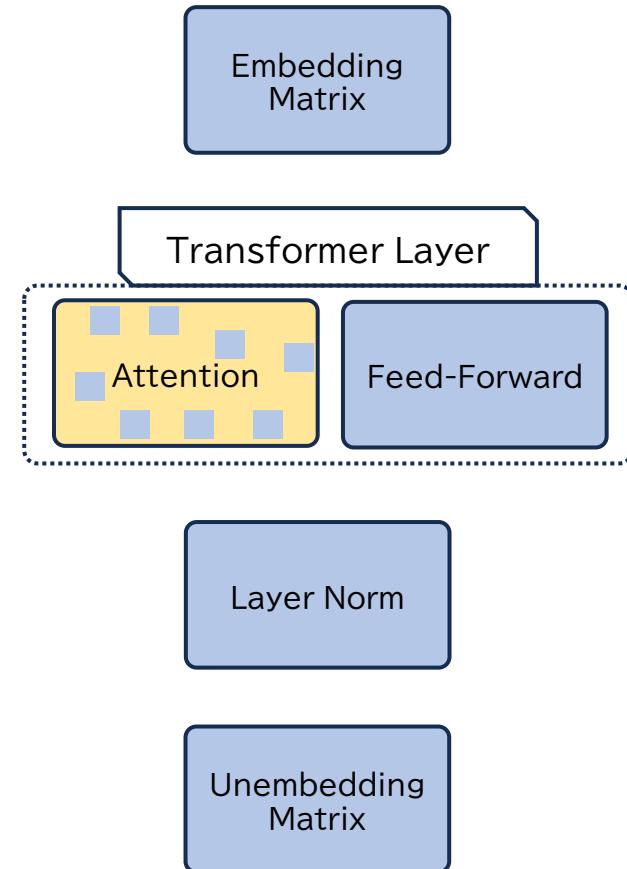
| Test Only                              |  | Train Only                         |   |
|--|--|------------------------------------|---|
| T-1: Missing-value identification (MV) | Identify the row and column position of the only missing cell in a given table         | T-10: List extraction (LE)         | Extract a structured table, from a list that lacks explicit column delimiters [9, 13, 19] |
| T-2: Column-finding (CF)               | Identify the column-name of a specific value that appears only once in a given table   | T-11: Head value matching (HVM)    | Match column-headers with its data values drawn from the same table                       |
| T-3: Table-QA (TQA)                    | Answer a natural-language question based on the content of a table ([11, 42, 49])      | T-12: Natural-language to SQL (NS) | Translate a natural-language question on a table into a SQL query ([62, 65])              |
| T-4: Column type annotation (CTA)      | Find the semantic type of a column, from a given list of choices ([16, 25, 63])        | T-13: Table summarization (TS)     | Produce a natural-language summary for the content in a table                             |
| Train / Test                           |  | T-14: Column augmentation (CA)     | Augment a table with additional columns compatible with a given table                     |
| T-5: Row-to-row transform (R2R)        | Transform table data based on input/output examples ([23, 24, 27])                     | T-15: Row augmentation (RA)        | Augment a table with additional rows compatible with a given table                        |
| T-6: Entity matching (EM)              | Match rows from two tables that refer to the same real-world entity ([32, 38, 41, 66]) | T-16: Row/column swapping (RCSW)   | Manipulate a given table, by swapping the position of two rows or columns                 |
| T-7: Schema matching (SM)              | Match columns from two tables that refer to the same meaning ([30, 36, 44])            | T-17: Row/column filtering (RCF)   | Manipulate a given table, by filtering on given rows or columns                           |
| T-8: Data imputation (DI)              | Predict the missing values in a cell based on the table context ([7, 37])              | T-18: Row/column sorting (RCS)     | Manipulate a given table, by performing sorting on given rows or columns                  |
| T-9: Error detection (ED)              | Detect data values in a table that is a likely error from misspelling ([14, 45])       |                                    |   |

# Table-GPT



# Parameter-Efficient Fine-Tuning (PEFT)

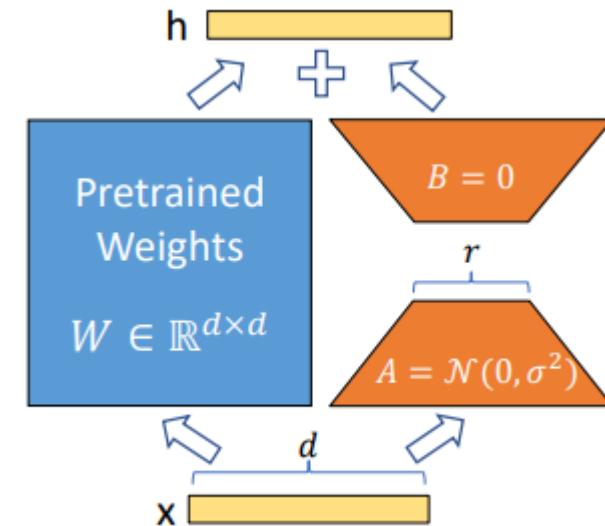
- Adjust only a **small portion** of parameters
- Introduce additional lightweight modules
- Advantages
  - Resource Efficiency
  - Flexibility
- Disadvantages
  - Limited Capacity
  - Complex Implementation
- Use Case
  - Task-specific data is limited
  - Computational resources are limited



# LoRA (Low-Rank Adaptation)

- Updates only low-rank matrices added to the model weights
- Advantages
  - Parameter Efficiency
  - Memory Efficiency
  - Modularity
  - Adaptability
- Disadvantages
  - Hyperparameter Optimization Needed

|                          | Weight Type          | $r = 1$ | $r = 2$ | $r = 4$ | $r = 8$ | $r = 64$ |
|--------------------------|----------------------|---------|---------|---------|---------|----------|
| WikiSQL( $\pm 0.5\%$ )   | $W_q$                | 68.8    | 69.6    | 70.5    | 70.4    | 70.0     |
|                          | $W_q, W_v$           | 73.4    | 73.3    | 73.7    | 73.8    | 73.5     |
|                          | $W_q, W_k, W_v, W_o$ | 74.1    | 73.7    | 74.0    | 74.0    | 73.9     |
| MultiNLI ( $\pm 0.1\%$ ) | $W_q$                | 90.7    | 90.9    | 91.1    | 90.7    | 90.7     |
|                          | $W_q, W_v$           | 91.3    | 91.4    | 91.3    | 91.6    | 91.4     |
|                          | $W_q, W_k, W_v, W_o$ | 91.2    | 91.7    | 91.7    | 91.5    | 91.4     |



$$h = W_0 x + \Delta W x = W_0 x + BAx$$

$$W = W_0 + BA$$

# SpreadSheetLLM

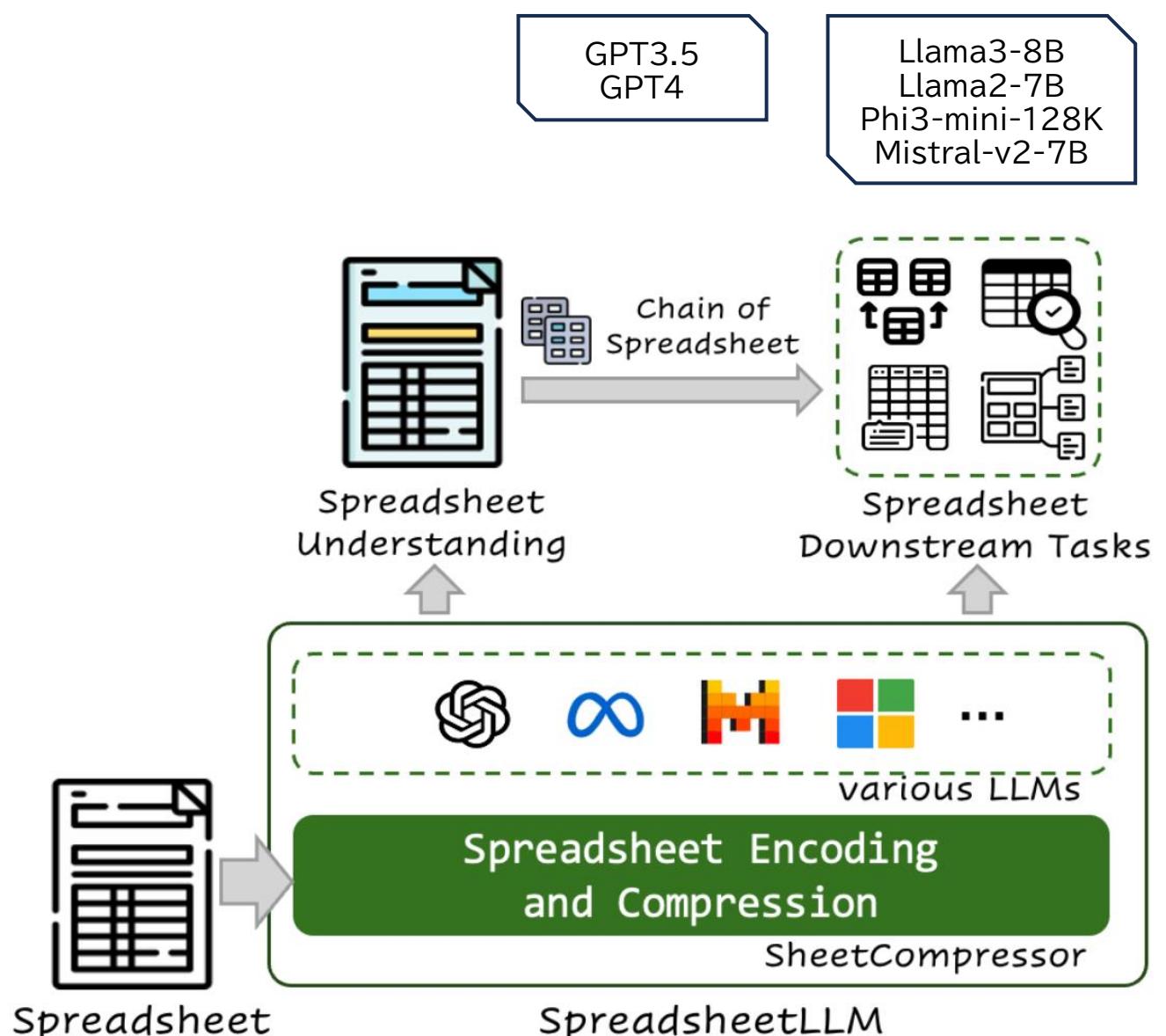
Table Understanding

Table Analysis

GPT3.5  
GPT4

Llama3-8B  
Llama2-7B  
Phi3-mini-128K  
Mistral-v2-7B

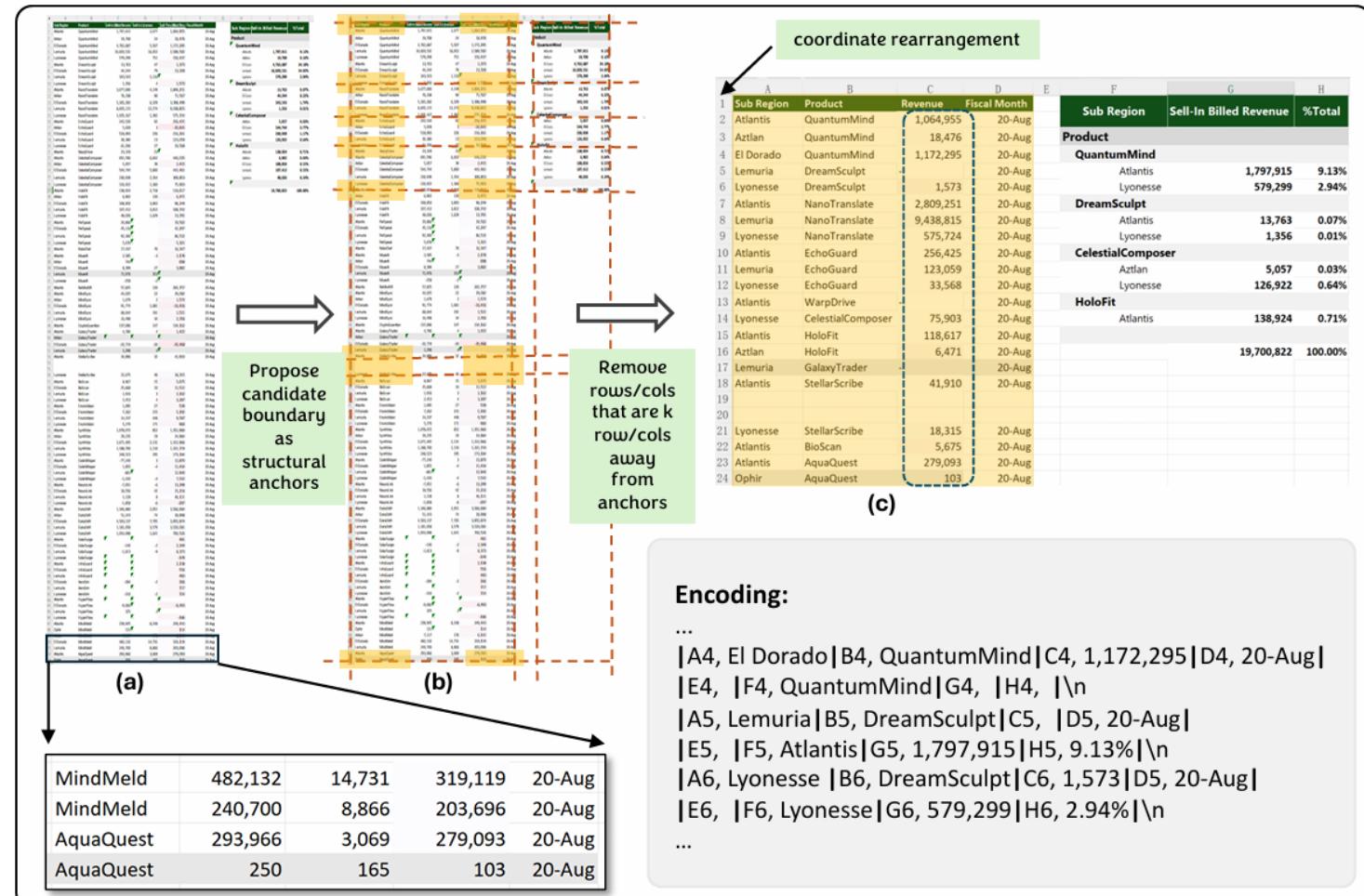
- Innovative Spreadsheet Encoding
- Chain of Spreadsheet Methodology
- Comprehensive Spreadsheet Analysis



# SpreadSheetLLM

- Structural-Anchor-Based Extraction

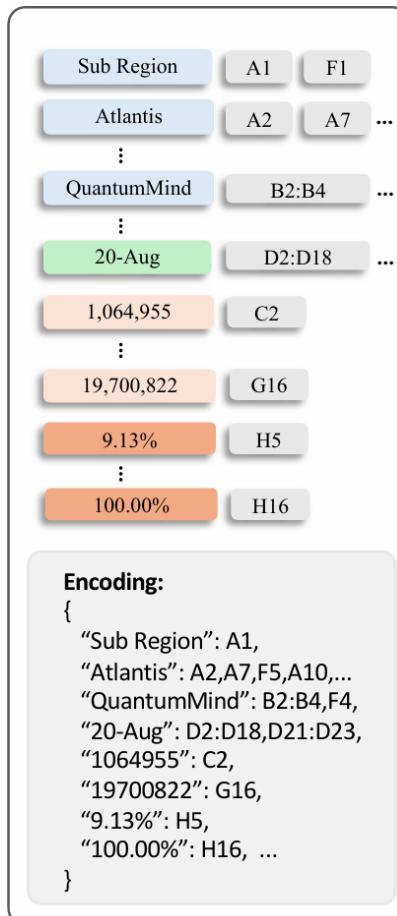
- Identifies key rows and columns at table boundaries
- Removes distant rows & columns to create a “skeleton” of spreadsheets



# SpreadSheetLLM

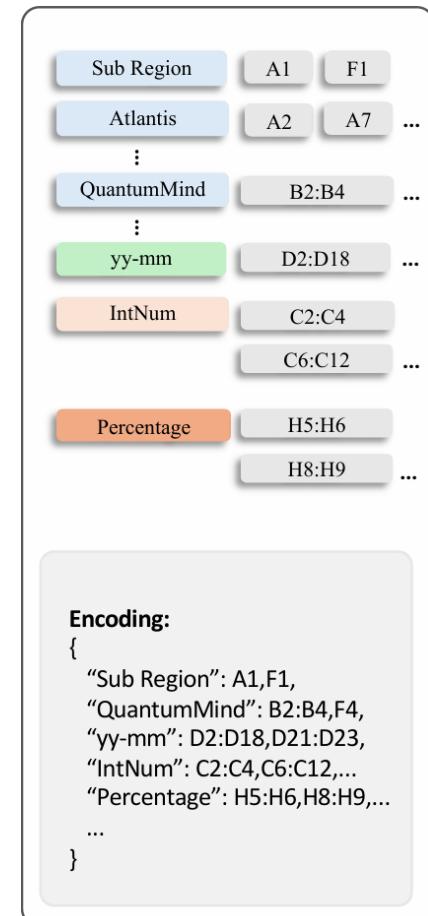
- Inverted-index Translation

- Translates spreadsheet data into a dictionary-like format
- Each unique value is stored once, with corresponding cell addresses indexed
- Empty cells are excluded from the encoding



- Data-format-aware Aggregation

- Groups cells into clusters
- Recognizes and aggregates formats



# SpreadSheetLLM

- Results: Table Detection

| Model & Method      | Small        | Medium       | Large        | Huge         | All          |
|---------------------|--------------|--------------|--------------|--------------|--------------|
| <b>ICL</b>          |              |              |              |              |              |
| Mistral-v2          | 0.071        | 0.013        | 0.029        | 0.017        | 0.036        |
| GPT4                | 0.318        | 0.292        | 0.090        | 0.000        | 0.154        |
| GPT4-compress       | 0.480        | 0.454        | 0.373        | 0.330        | 0.410        |
| <b>Fine-tune</b>    |              |              |              |              |              |
| Llama3              | 0.715        | 0.765        | 0.290        | 0.000        | 0.471        |
| Llama2              | 0.557        | 0.378        | 0.107        | 0.000        | 0.280        |
| Phi3                | 0.604        | 0.481        | 0.201        | 0.130        | 0.330        |
| Mistral-v2          | 0.700        | 0.784        | 0.472        | 0.123        | 0.542        |
| GPT4                | 0.779        | 0.707        | 0.288        | 0.000        | 0.520        |
| Llama3-compress     | 0.825        | 0.768        | 0.664        | 0.617        | 0.719        |
| Llama2-compress     | 0.710        | 0.722        | 0.633        | 0.578        | 0.660        |
| Phi3-compress       | 0.800        | 0.673        | 0.624        | 0.675        | 0.689        |
| Mistral-v2-compress | 0.778        | 0.729        | 0.686        | 0.744        | 0.726        |
| GPT3.5-compress     | 0.795        | 0.649        | 0.600        | 0.680        | 0.680        |
| GPT4-compress       | 0.810        | <b>0.832</b> | 0.718        | 0.690        | 0.759        |
| -w/o Aggregation    | <b>0.864</b> | 0.816        | <b>0.739</b> | <b>0.753</b> | <b>0.789</b> |
| TableSense-CNN      | 0.785        | 0.788        | 0.567        | 0.561        | 0.666        |

- Results: Spreadsheet QA

| Model                          | Accuracy     |
|--------------------------------|--------------|
| TAPEX                          | 0.378        |
| Binder                         | 0.622        |
| GPT4                           | 0.466        |
| GPT4-compress-w/o splitting    | 0.651        |
| GPT4-compress-w/o splitting-FT | 0.694        |
| GPT4-compress                  | 0.684        |
| GPT4-compress-FT               | <b>0.743</b> |

# TAT-LLM

## Table Analysis

- Hybrid Tabular and Textual QA
- Step-wise Pipeline
- External Executor

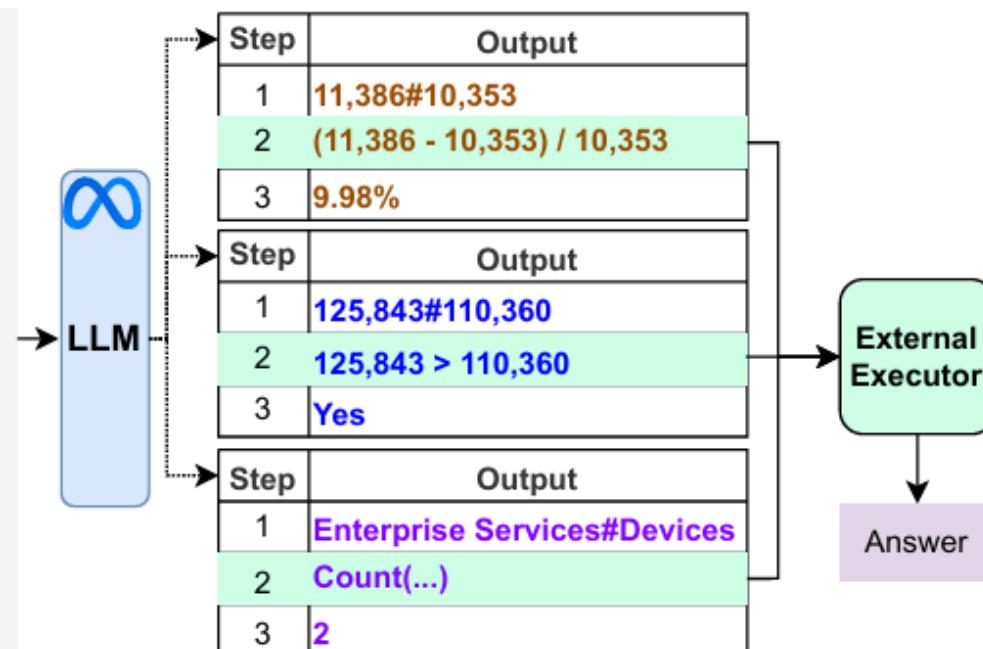
Llama2-7B-Chat

**Table** ### Instruction  
→ Please complete the task in three steps:  
1. **Extractor**: extract the relevant values ...  
2. **Reasoner**: generate the logic or equation ...  
3. **Executor**: calculate the answer ...  
Please organize the results in the following table:  
| Step | Output |

**Q1** → | 1 | |  
| 2 | |  
| 3 | |

**Q2** → ### Table  
→ {Table} in markdown format  
### Text  
{Texts}

**Q3** → ### Question  
{Question}  
### Response



# TAT-LLM

---

**Algorithm 1 :** External Executor
 

---

**Input**  $O_1$ : the output of *Extractor*;  $O_2$ : the output of *Reasoner*;  $O_3$ : the output of *Executor*;  $Q_t$ : the predicted question type.

```

1:  $answer \leftarrow O_3$ 
2: if  $O_2$  is a valid arithmetic equation then
3:    $answer \leftarrow round(eval(O_2), 4)$ 
4: else if “#” in  $O_2$  then # multiple values are
   separated with “#”
5:    $arr \leftarrow O_2.split("#")$ 
6:    $answer \leftarrow len(arr)$ 
7: else if “>” in  $O_2$  or “<” in  $O_2$  then
8:    $answer \leftarrow eval(O_2)$ 
9: else if  $O_2 = "N.A."$  then
10:  if  $Q_t = "Span"$  then
11:     $answer \leftarrow O_1$ 
12:  else if  $Q_t = "Multiple Spans"$  then
13:     $arr \leftarrow O_1.split("#")$  # spans are sep-
       arated with “#”
14:     $answer \leftarrow arr$ 
15:  end if
16: end if
```

---

| Model                        | FinQA |       | TAT-QA         |       | TAT-DQA        |  |
|------------------------------|-------|-------|----------------|-------|----------------|--|
|                              | EM    | EM    | F <sub>1</sub> | EM    | F <sub>1</sub> |  |
| GPT-4                        | 63.91 | 71.92 | 79.71          | 64.46 | 72.20          |  |
| TAT-LLM <sub>All</sub> (7B)  |       |       |                |       |                |  |
| w/o <i>External Executor</i> | 48.47 | 58.69 | 67.21          | 54.84 | 63.68          |  |
| w <i>External Executor</i>   | 65.13 | 76.49 | 85.13          | 71.38 | 80.24          |  |
| gains (+)                    | 16.66 | 17.80 | 17.92          | 16.54 | 16.56          |  |
| TAT-LLM <sub>All</sub> (13B) |       |       |                |       |                |  |
| w/o <i>External Executor</i> | 60.05 | 62.60 | 70.73          | 59.95 | 68.61          |  |
| w <i>External Executor</i>   | 71.75 | 76.79 | 85.05          | 71.86 | 80.50          |  |
| gains (+)                    | 11.70 | 14.19 | 14.32          | 11.91 | 11.89          |  |
| TAT-LLM <sub>All</sub> (70B) |       |       |                |       |                |  |
| w/o <i>External Executor</i> | 70.10 | 76.61 | 83.55          | 71.74 | 78.99          |  |
| w <i>External Executor</i>   | 76.81 | 81.42 | 88.49          | 76.55 | 83.90          |  |
| gains (+)                    | 6.71  | 4.81  | 4.94           | 4.81  | 4.91           |  |

# TAT-LLM

| Type                            | Model                       | EM                               |
|---------------------------------|-----------------------------|----------------------------------|
| <b>Human Expert Performance</b> |                             | 91.16                            |
| <b>Fine-tuned</b>               | Longformer                  | 21.90                            |
|                                 | NeRd                        | 48.57                            |
|                                 | FinQANet <sub>BERT</sub>    | 50.00                            |
|                                 | DyRRen <sub>BERT</sub>      | 59.37                            |
|                                 | FinQANet <sub>RoBERTa</sub> | 61.24                            |
|                                 | ELASTIC <sub>RoBERTa</sub>  | 62.66                            |
|                                 | DyRRen <sub>RoBERTa</sub>   | 63.30                            |
| <b>Zero-shot</b>                | Vicuna (7B)                 | 10.11                            |
|                                 | LLaMA 2-Chat (7B)           | 15.43                            |
|                                 | LLaMA 2-Chat (70B)          | 32.17                            |
|                                 | MAmmoTH (70B)               | 36.09                            |
|                                 | WizardMath (70B)            | 47.25                            |
|                                 | GPT3.5-Turbo                | 58.00                            |
|                                 | GPT-4                       | <u>63.91</u>                     |
| <b>Ours</b>                     | TAT-LLM (7B)                | ( <b>+0.69</b> )<br><b>64.60</b> |

Table 2: Performance of our TAT-LLM model and compared models on the test set of FinQA. Best results are marked in bold and numbers in red indicate the improvement over the underlined second-best results.

| Type                            | Model              | EM           | F <sub>1</sub>                   |
|---------------------------------|--------------------|--------------|----------------------------------|
| <b>Human Expert Performance</b> |                    | 84.1         | 90.8                             |
| <b>Fine-tuned</b>               | TagOp              | 50.10        | 58.00                            |
|                                 | TeaBReaC           | 55.80        | 63.80                            |
|                                 | KIQA               | 58.20        | 67.40                            |
|                                 | FinMath            | 58.30        | 68.20                            |
|                                 | GANO               | 61.90        | 72.10                            |
|                                 | MHST               | 63.60        | 72.70                            |
|                                 | UniPCQA            | 63.90        | 72.20                            |
|                                 | SoarGraph          | 65.40        | 75.30                            |
|                                 | UniRPG             | 67.20        | 76.00                            |
|                                 | RegHNT             | 70.30        | 77.90                            |
|                                 | MVGE               | 70.90        | 79.10                            |
|                                 | Vicuna (7B)        | 32.53        | 40.97                            |
| <b>Zero-shot</b>                | LLaMA 2-Chat (7B)  | 37.16        | 45.37                            |
|                                 | MAmmoTH (70B)      | 38.97        | 46.51                            |
|                                 | WizardMath (70B)   | 39.63        | 45.28                            |
|                                 | LLaMA 2-Chat (70B) | 45.94        | 53.80                            |
|                                 | GPT3.5-Turbo       | 59.47        | 68.11                            |
|                                 | GPT-4              | <u>71.92</u> | <u>79.71</u>                     |
|                                 | <b>Ours</b>        | TAT-LLM (7B) | ( <b>+2.64</b> )<br><b>74.56</b> |
|                                 |                    |              | ( <b>+3.17</b> )<br><b>82.88</b> |

Table 3: Performance of our TAT-LLM model and compared models on the test set of TAT-QA.

| Type                            | Model              | EM                               | F <sub>1</sub>                   |
|---------------------------------|--------------------|----------------------------------|----------------------------------|
| <b>Human Expert Performance</b> |                    | 84.1                             | 90.8                             |
| <b>Fine-tuned</b>               | NumNet+ V2         | 30.60                            | 40.10                            |
|                                 | TagOp              | 33.70                            | 42.50                            |
|                                 | MHST               | 41.50                            | 50.70                            |
|                                 | Doc2SoarGraph      | 59.20                            | 67.60                            |
| <b>Zero-shot</b>                | Vicuna (7B)        | 28.44                            | 36.72                            |
|                                 | LLaMA 2-Chat (7B)  | 34.52                            | 42.32                            |
|                                 | MAmmoTH (70B)      | 35.42                            | 42.82                            |
|                                 | WizardMath (70B)   | 36.44                            | 41.55                            |
|                                 | LLaMA 2-Chat (70B) | 41.91                            | 49.74                            |
|                                 | GPT3.5-Turbo       | 52.74                            | 61.40                            |
|                                 | <b>GPT-4</b>       | <u>64.46</u>                     | <u>72.20</u>                     |
| <b>Ours</b>                     | TAT-LLM (7B)       | ( <b>+4.99</b> )<br><b>69.45</b> | ( <b>+5.55</b> )<br><b>77.75</b> |

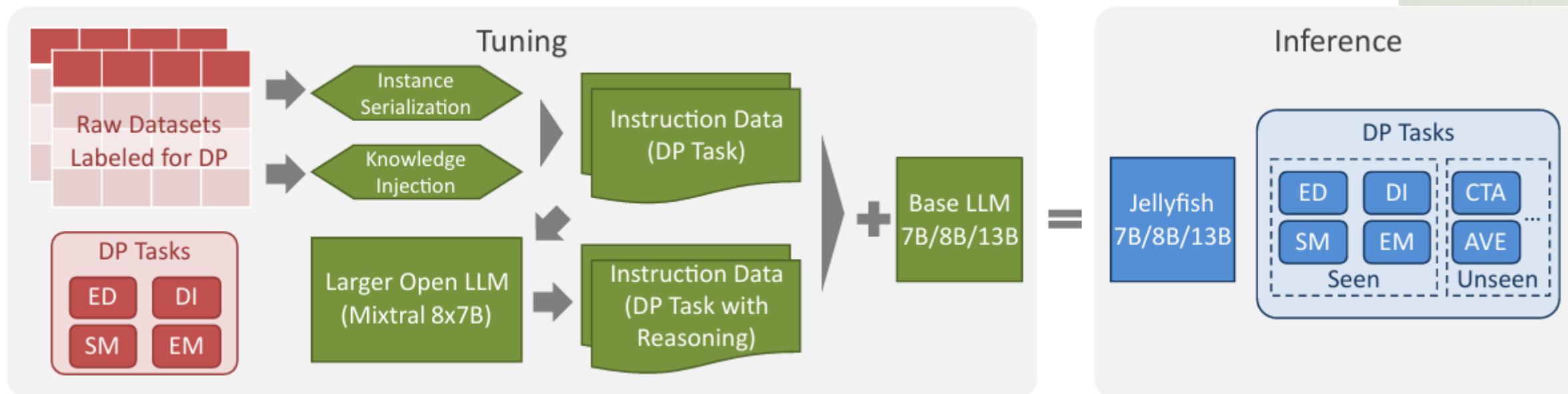
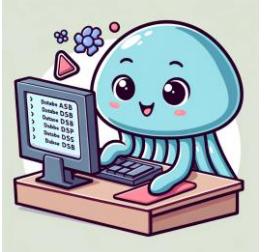
Table 4: Performance of our TAT-LLM model and compared models on the test set of TAT-DQA.

# Jellyfish

## Table Preprocessing

- Specialized in Data Preprocessing
- Data Preprocessing Interpreter
- Jellyfish Reasoning Dataset

OpenOrca-Platypus2-13B  
Llama-3.1-8B-Instruct  
Mistral-7B-Instruct-v0.2



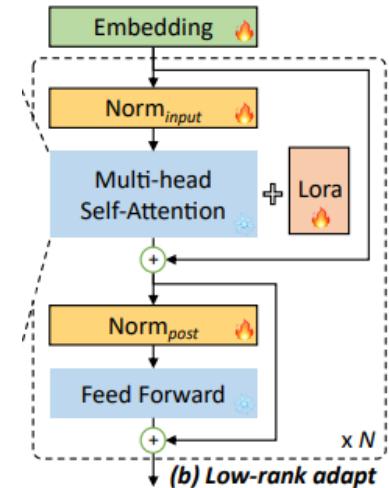
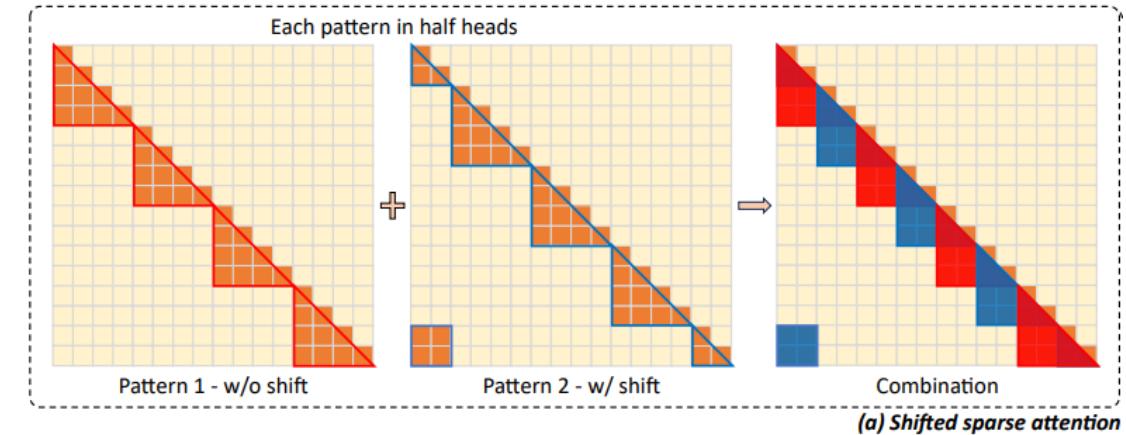
# Jellyfish

| Task | Type    | Dataset   | Model   |   |   |   |  |   |  |  |  |
|------|---------|---|---|---|---|---|--|---|--|--|--|
|      |         |   | Best of non-LLM   | GPT-3   | GPT-3.5   | GPT-4   | GPT-4o   | Table-GPT   | Jellyfish-7B   | Jellyfish-8B   | Jellyfish-13B  |
| ED   | Seen    | Adult Hospital  | 99.10<br>94.40  | 99.10<br><b>97.80</b>   | 92.01<br>90.74  | 92.01<br>90.74  | 83.58<br>44.76                                     | –   | 77.40<br>94.51   | 73.74<br>93.40   | <b>99.33</b><br><u>95.59</u>                                   |
|      | Unseen  | Flights Rayyan  | 81.00<br>79.00  | –<br>–  | –<br>–  | <b>83.48</b><br><u>81.95</u>                                      | 66.01<br>68.53                                     | –   | 69.15<br>75.07   | 66.21<br>81.06   | <u>82.52</u><br><b>90.65</b>                                   |
| DI   | Seen    | Buy Restaurant  | 96.50<br>77.20  | 98.50<br>88.40  | 98.46<br><u>94.19</u>                                   | <b>100</b><br><b>97.67</b>  | <b>100</b><br>90.70                                | –   | 98.46<br>89.53   | 98.46<br>87.21   | <b>100</b><br>89.53  |
|      | Unseen  | Flipkart Phone  | 68.00<br>86.70  | –<br>–  | –<br>–  | <b>89.94</b><br><b>90.79</b>                                      | 83.20<br>86.78                                     | –   | 87.14<br>86.52   | <u>87.48</u><br>85.68  | 81.68<br><u>87.21</u>  |
| SM   | Seen    | MIMIC-III Synthea   | 20.00<br>38.50  | –<br>45.20  | –<br><u>57.14</u>                                       | 40.00<br><b>66.67</b>   | 29.41<br>6.56                                      | –   | <b>53.33</b><br>55.56  | <u>45.45</u><br>47.06  | 40.00<br>56.00   |
|      | Unseen  | CMS   | <b>50.00</b>  | –   | –   | 19.35   | 22.22  | –   | 42.86  | 38.10  | <b>59.29</b>   |
|      | Seen    | Amazon-Google<br>Beer<br>DBLP-ACM<br>DBLP-GoogleScholar<br>Fodors-Zagats<br>iTunes-Amazon | 75.58<br>94.37<br><b>98.99</b><br><u>95.70</u><br><b>100</b><br>97.06 | 63.50<br><b>100</b><br>96.60<br>83.80<br><b>100</b><br><u>98.20</u> | 66.50<br>96.30<br>96.99<br>76.12<br><b>100</b><br>96.40 | 74.21<br><b>100</b><br>97.44<br>91.87<br><b>100</b><br><b>100</b> | 70.91<br>90.32<br>95.87<br>90.45<br>93.62<br>98.18 | 70.10<br>96.30<br>93.80<br>92.40<br><b>100</b><br>94.30 | <b>81.69</b><br><b>100.00</b><br>98.65<br>94.88<br><b>100</b><br>96.30 | <u>81.42</u><br><b>100.00</b><br>98.77<br>95.03<br><b>100</b><br>96.30 | 81.34<br>96.77<br><u>98.98</u><br><b>98.51</b><br>100<br>98.11 |
| EM   | Unseen  | Abt-Buy<br>Walmart-Amazon   | 89.33<br>86.89  | –<br>87.00  | –<br>86.17  | <b>92.77</b><br><b>90.27</b>                                      | 78.73<br>79.19                                     | –<br>82.40  | 86.06<br>84.91   | 88.84<br>85.24   | <u>89.58</u><br><u>89.42</u>                                   |
|      | Average |   | 80.44   | -   | -   | 84.17   | 72.58  | -   | 82.74  | 81.55  | <b>86.02</b>   |

| Task | Dataset            | Model               |                     |                     |                |                       |                       |                |                |                       |                |
|------|--------------------|---------------------|---------------------|---------------------|----------------|-----------------------|-----------------------|----------------|----------------|-----------------------|----------------|
|      |                    | RoBERTa (159 shots) | RoBERTa (356 shots) | Stable Beluga 2 70B | SOLAR 70B      | GPT-3.5               | GPT-4                 | GPT-4o         | Jellyfish-7B   | Jellyfish-8B          | Jellyfish-13B  |
| CTA  | SOTAB              | 79.20               | 89.73               | –                   | –              | 89.47                 | <b>91.55</b>          | 65.06          | 83.00          | 76.33                 | 82.00          |
| AVE  | AE-110k<br>OA-Mine | –                   | –                   | 52.10<br>50.80      | 49.20<br>55.20 | <b>61.30</b><br>62.70 | 55.50<br><b>68.90</b> | 55.77<br>60.20 | 56.09<br>51.98 | <u>59.55</u><br>59.22 | 58.12<br>55.96 |

# LongLoRA

- Extends the context sizes of large language models with limited computational cost
- Shifted Sparse Attention (S2-Attn)
- Advantages
  - Efficient Context Extension
  - Compatible with existing optimizations
- Disadvantages
  - Requires careful handling of attention mechanisms



# TableLlama

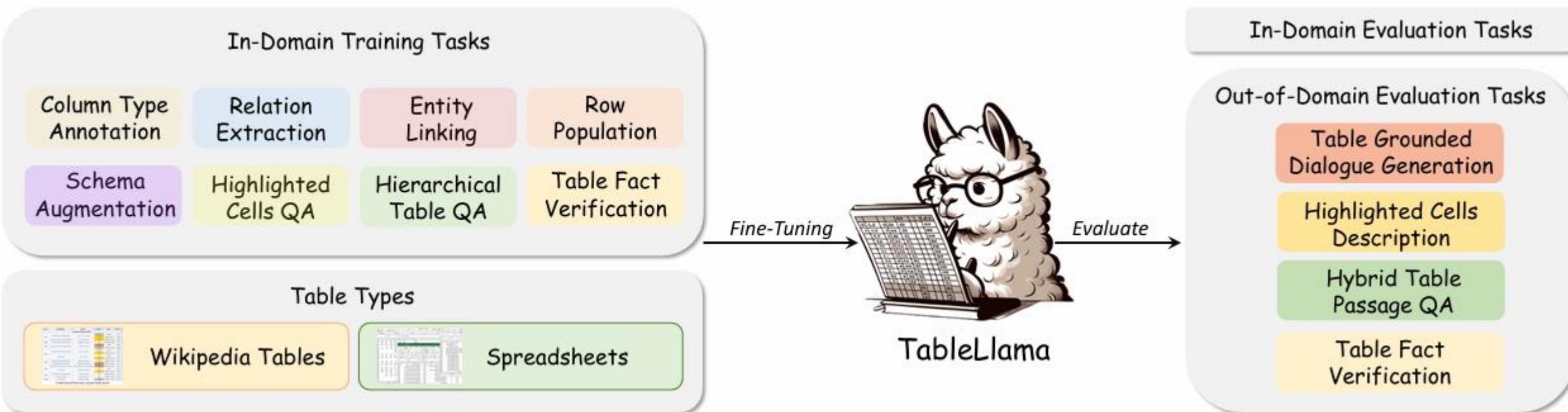
Table Preprocessing

Table Understanding

Table Analysis

Llama2-7B

- TableInstruct Dataset
- Real-world data sources
- LongLoRA for Extending Context Window



# TableLlama

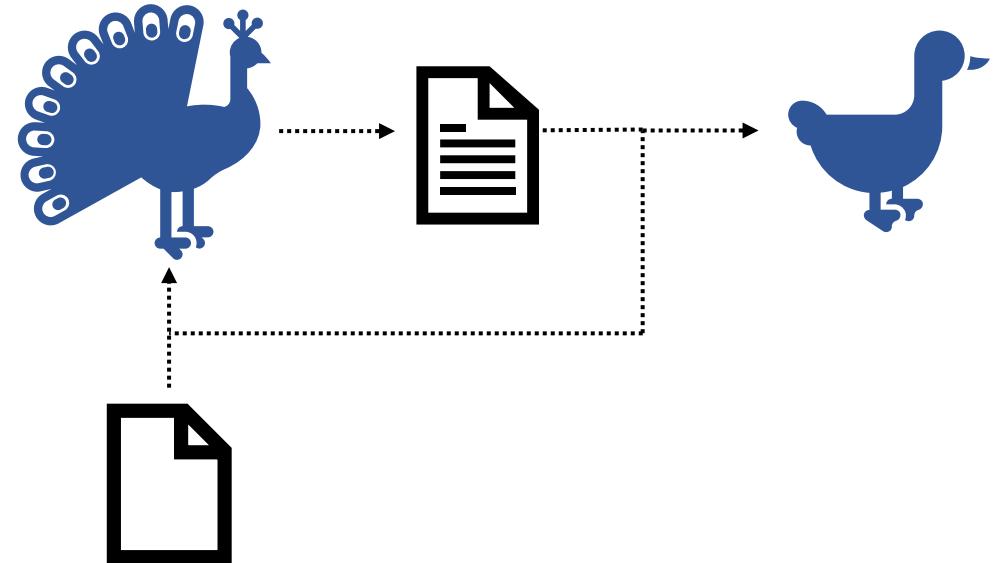
| Task Category        | Task Name                          | Dataset                          | In-domain | #Train<br>(Table/Sample) | #Test<br>(Table/Sample) | Input<br>min | Token<br>max | Length<br>median |
|----------------------|------------------------------------|----------------------------------|-----------|--------------------------|-------------------------|--------------|--------------|------------------|
| Table Interpretation | Col Type Annot.                    | TURL (Deng et al., 2020)         | Yes       | 397K/628K                | 1K/2K                   | 106          | 8192         | 2613             |
|                      | Relation Extract.                  |                                  | Yes       | 53K/63K                  | 1K/2K                   | 2602         | 8192         | 3219             |
|                      | Entity Linking                     |                                  | Yes       | 193K/1264K               | 1K/2K                   | 299          | 8192         | 4667             |
| Table Augmentation   | Schema Aug.                        | TURL (Deng et al., 2020)         | Yes       | 288K/288K                | 4K/4K                   | 160          | 1188         | 215              |
|                      | Row Pop.                           |                                  | Yes       | 286K/286K                | 0.3K/0.3K               | 264          | 8192         | 1508             |
| Question Answering   | Hierarchical Table QA              | HiTab (Cheng et al., 2022b)      | Yes       | 3K/7K                    | 1K/1K                   | 206          | 5616         | 978              |
|                      | Highlighted Cells QA               | FeTaQA (Nan et al., 2022)        | Yes       | 7K/7K                    | 2K/2K                   | 261          | 5923         | 740              |
|                      | Hybrid Table QA                    | HybridQA (Chen et al., 2020b)    | No        | –                        | 3K/3K                   | 248          | 2497         | 675              |
|                      | Table QA                           | WikiSQL (Zhong et al., 2017)     | No        | –                        | 5K/16K                  | 198          | 2091         | 575              |
|                      | Table QA                           | WikiTQ (Pasupat and Liang, 2015) | No        | –                        | 0.4K/4K                 | 263          | 2688         | 709              |
| Fact Verification    | Fact Verification                  | TabFact (Chen et al., 2020a)     | Yes       | 16K/92K                  | 2K/12K                  | 253          | 4975         | 630              |
|                      |                                    | FEVEROUS (Aly et al., 2021)      | No        | –                        | 4K/7K                   | 247          | 8192         | 648              |
| Dialogue Generation  | Table Grounded Dialogue Generation | KVRET (Eric et al., 2017)        | No        | –                        | 0.3K/0.8K               | 187          | 1103         | 527              |
| Data-to-Text         | Highlighted Cells Description      | ToTTo (Parikh et al., 2020)      | No        | –                        | 7K/8K                   | 152          | 8192         | 246              |

# TableLlama

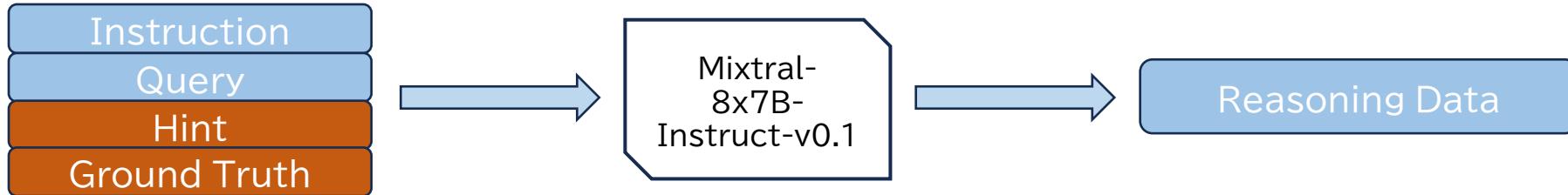
| In-domain Evaluation     |          |       |              |                                     |                 |         |        |
|--------------------------|----------|-------|--------------|-------------------------------------|-----------------|---------|--------|
| Datasets                 | Metric   | Base  | TableLlama   | SOTA                                | GPT-3.5         | GPT-4§  |        |
| Column Type Annotation   | F1       | 3.01  | 94.39        | <b>94.54*</b> † (Deng et al., 2020) | 30.88           | 31.75   |        |
| Relation Extraction      | F1       | 0.96  | 91.95        | <b>94.91*</b> † (Deng et al., 2020) | 27.42           | 52.95   |        |
| Entity Linking           | Accuracy | 31.80 | <b>93.65</b> | 84.90*† (Deng et al., 2020)         | 72.15           | 90.80   |        |
| Schema Augmentation      | MAP      | 36.75 | <b>80.50</b> | 77.55*† (Deng et al., 2020)         | 49.11           | 58.19   |        |
| Row Population           | MAP      | 4.53  | 58.44        | <b>73.31*</b> † (Deng et al., 2020) | 22.36           | 53.40   |        |
| HiTab                    | Exec Acc | 14.96 | <b>64.71</b> | 47.00*† (Cheng et al., 2022a)       | 43.62           | 48.40   |        |
| FeTaQA                   | BLEU     | 8.54  | <b>39.05</b> | 33.44 (Xie et al., 2022)            | 26.49           | 21.70   |        |
| TabFact                  | Accuracy | 41.65 | 82.55        | <b>84.87*</b> (Zhao and Yang, 2022) | 67.41           | 74.40   |        |
| Out-of-domain Evaluation |          |       |              |                                     |                 |         |        |
| Datasets                 | Metric   | Base  | TableLlama   | SOTA                                | $\Delta_{Base}$ | GPT-3.5 | GPT-4§ |
| FEVEROUS                 | Accuracy | 29.68 | 73.77        | 85.60 (Tay et al., 2022)            | +44.09          | 60.79   | 71.60  |
| HybridQA                 | Accuracy | 23.46 | 39.38        | 65.40* (Lee et al., 2023)           | +15.92          | 40.22   | 58.60  |
| KVRET                    | Micro F1 | 38.90 | 48.73        | 67.80 (Xie et al., 2022)            | +9.83           | 54.56   | 56.46  |
| ToTTo                    | BLEU     | 10.39 | 20.77        | 48.95 (Xie et al., 2022)            | +10.38          | 16.81   | 12.21  |
| WikiSQL                  | Accuracy | 15.56 | 50.48        | 92.70 (Xu et al., 2023b)            | +34.92          | 41.91   | 47.60  |
| WikiTQ                   | Accuracy | 29.26 | 35.01        | 57.50† (Liu et al., 2022)           | +5.75           | 53.13   | 68.40  |

# Knowledge Distillation

- Training a **student model** using the output of a **teacher model**
- Advantages
  - Model Compression
  - Enhanced Generalization
- Disadvantages
  - Rely on the teacher model
  - Potential Performance Trade-off
- Use Case
  - Lake of labeled training data
  - Smaller model is required



# Jellyfish



**(system message)** You are an AI assistant that follows instruction extremely well. User will give you a question. Your task is to answer as faithfully as you can. While answering, provide detailed explanation and justify your answer.

**(task description)** You are tasked with determining whether two Products listed below are the same based on the information provided. Carefully compare all the attributes before making your decision.

**(injected knowledge)** Note that missing values (N/A or "nan") should not be used as a basis for your decision. Note that different factories can belong to the same parent company. The company name of Product B may occur in its product name.

**(instance content)** Product A: [name: "Sequoia American Amber Ale", factory: "Wig And Pen"] Product B: [name: "Aarhus Cains Triple A American Amber Ale", factory: "Aarhus Bryghus"]

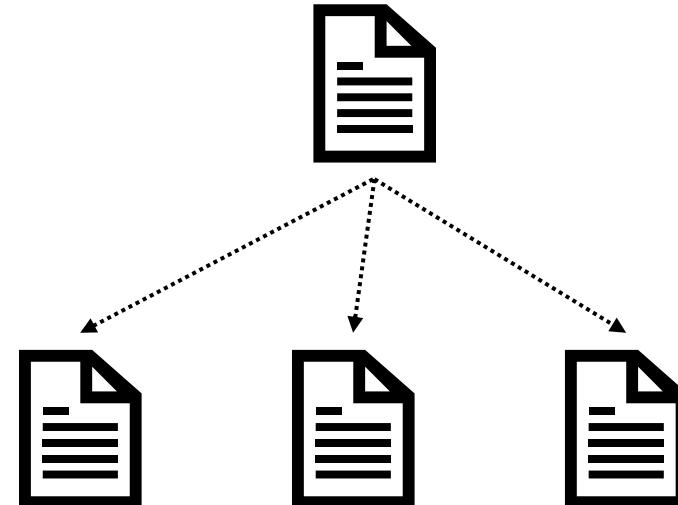
**(question)** Are Product A and Product B the same Product?

**(output format)** After your reasoning, finish your response in a separate line with and ONLY with your final answer. Choose your final answer from [Yes, No].

**(answer hint)** You can use the "Hint" below, but your response cannot contain any information from it. Hint: the final answer is "No"

# Data Augmentation Enhanced Fine-Tuning

- Incorporates data augmentation to enrich the training dataset
- Advantages
  - Improved Generalization
  - Enhanced Robustness
  - Data Efficiency
- Disadvantages
  - Potential for Mislabeling
  - Diminishing Returns
- Use Case
  - Limited Data Scenarios
  - Overfitting Issues



# Table-GPT

- Synthesize Diverse Table-Tasks

- Task Synthesis for Task-Diversity
  - Creating new, diverse tasks from real-world table data
- Data-Diversity through Test Case Synthesis
  - Generate multiple test cases for a broader range of table scenarios

---

**Algorithm 1:** Synthesize table-tasks for table-tuning

---

```
input  : A corpus of diverse real tables C, a set of table-task types S
output : Diverse synthesized table-tasks A = { (Ins, T, C) }
1  D ← {}, A ← {}
2  foreach T ∈ C, S ∈ S do
3      | (Ins, T, C) ← Synthesize-Table-Task(S, T) // (Section 4.2)
4      | D ← D ∪ (Ins, T, C)
5  foreach (Ins, T, C) ∈ D do
6      | Ins' ← Augment-Instruction(Ins) // (Section 4.3)
7      | T' ← Augment-Table(T) // (Section 4.3)
8      | C' ← Augment-Completion(C) // (Section 4.3)
9      | A ← A ∪ (Ins', T', C')
10 return A
```

---

# Table-GPT

- Augment Synthesized Table-Tasks

- Instruction-Level Augmentation
- Table-Level Augmentation
- Completion-Level Augmentation

## T-8: Entity Matching

**Instruction:** Please determine whether the first and second row below may refer to the same entity. Please answer "yes/no" in JSON.

**Row-1:**

|       |           |       |             |     |
|-------|-----------|-------|-------------|-----|
| Brand | Model     | Price | Description | ... |
| Apple | iPhone-15 | \$995 | ...         | ... |

**Row-2:**

|       |            |        |             |     |
|-------|------------|--------|-------------|-----|
| Brand | Model      | Price  | Description | ... |
| Apple | iPhone pro | \$1500 | ...         | ... |

**Model response:**

The first and second product are likely different, because they have very different prices.  
Therefore, the answer is {"answer": "No"}.

## T-9: Error Detection (misspellings)

**Instruction:** Please examine the table below and check which cell may be misspelled. If there is no erroneous cell, return 'None'. Please answer using JSON: {"answer": ""}.

**Table:**

|                  |         |             |     |
|------------------|---------|-------------|-----|
| State            | Capital | Population  | ... |
| Washington       | Olympia | 7.7 million | ... |
| <u>Missisipi</u> | Jackson | 2.9 million | ... |

**Model response:**

The cell "Missisipi" is misspelled, and should be corrected as " Mississippi". Therefore, the answer is {"answer": "Missisipi".}

# Table Data Representation & Prompt Design

# Table-GPT

## Header Value Matching Zero-Shot

### Prompt:

Task: Given the input table data and the list of candidate headers, please determine the most suitable column header for each column in the table. Please only choose column headers from the candidate list. Please only return the most suitable column header for each column. Return the chosen column headers in a list. Do not return the entire table. Return the final result as JSON in the format {"column\_headers": "<a list of headers for each column chosen from the candidate list>"}.

[Q]:

\*\*Table Data:\*\*

|||

|     |     |           |
|-----|-----|-----------|
| --- | --- | ---       |
| 1   | 681 | Multihull |
| 3   | 911 | Keelboat  |
| 2   | 947 | Dinghy    |
| nan | 920 | Dinghy    |
| 1   | 870 | Dinghy    |

\*\*Candidate column headers:\*\*

- Portsmouth Number
- crew
- Type

Return the final result as JSON in the format {"column\_headers": "<a list of headers for each column chosen from the candidate list>"}.

[A]:

### Completion:

```
{"column_headers": ["crew", "Portsmouth Number", "Type"]}
```

# SpreadSheetLLM

## Prompt Template for Spreadsheet Table Detection

**(INSTRUCTION:)** Given an input that is a string denoting data of cells in an Excel spreadsheet. The input spreadsheet contains many tuples, describing the cells with content in the spreadsheet. Each tuple consists of two elements separated by a '|': the cell content and the cell address/region, like (Year|A1), ( |A1) or(IntNum|A1:B3). The content in some cells such as '#,##0'/'d-mmm-yy'/'H:mm:ss', etc., represents the CELL DATA FORMATS of Excel. The content in some cells such as 'IntNum'/'DateData'/'EmailData', etc., represents a category of data with the same format and similar semantics. For example, 'IntNum' represents integer type data, and 'ScientificNum' represents scientific notation type data. 'A1:B3' represents a region in a spreadsheet, from the first row to the third row and from column A to column B. Some cells with empty content in the spreadsheet are not entered. Now you should tell me the range of the table in a format like A2:D5, and the range of the table should only CONTAIN HEADER REGION and the data region. DON'T include the title or comments. Note that there can be more than one table in a string, so you should return all the RANGE. DON'T ADD OTHER WORDS OR EXPLANATION.

**(INPUT)** [Encoded Spreadsheet]

# TAT-LLM

## Template for FinQA dataset following Step-wise Pipeline.

Below is an instruction that describes a question answering task in the finance domain, paired with an input table and its relevant text that provide further context. The given question is relevant to the table and text. Generate an appropriate answer to the given question.

**(Instruction)** Given a table and a list of texts in the following, what is the answer to the question? Please complete the task in three steps:

**(Extractor)** In the first step, extract the relevant numerical values from the provided table or texts. Store these in the variable '{evidence}'. If there are multiple values, separate them using the '#' symbol.

**(Reasoner)** In the second step, generate an equation using the extracted numerical values. Store this equation in the variable '{equation}'.

**(Executor:)** In the third step, calculate the answer based on the equation and store it in the variable '{answer}'. Please organize the results in the following table:

| step | output     |
|------|------------|
| 1    | {evidence} |
| 2    | {equation} |
| 3    | {answer}   |

Finally, present the calculated answer in the format: "The answer is: {answer}"

**(Table) (Text) (Question)**

# Jellyfish

## Direct and reasoning prompt Template for entity matching

|                    | DP Task Data  | DP Task with Reasoning Data   |
|--------------------|---|---|
| system message     | You are an AI assistant that follows instruction extremely well. User will give you a question. Your task is to answer as faithfully as you can.                                    | You are an AI assistant that follows instruction extremely well. User will give you a question. Your task is to answer as faithfully as you can. While answering, provide detailed explanation and justify your answer. |
| task description   | You are tasked with determining whether two Products listed below are the same based on the information provided. Carefully compare all the attributes before making your decision. |   |
| injected knowledge | Note that missing values (N/A or "nan") should not be used as a basis for your decision.  |   |
| instance content   | Product A: [name: "Sequoia American Amber Ale", factory: "Wig And Pen"]<br>Product B: [name: "Aarhus Cains Triple A American Amber Ale", factory: "Aarhus Bryghus"]                 |   |
| question           | Are Product A and Product B the same Product?   |   |
| output format      | Choose your answer from: [Yes, No]  | After your reasoning, finish your response in a separate line with and ONLY with your final answer. Choose your final answer from [Yes, No].  |

# TableLlama

## ### Instruction:

This is a table **row population** task. The goal of this task is to populate the possible entities of the selected column for a table, given the Wikipedia page title, ... You will be given a list of entity candidates. Please rank them so that the most likely entities come first.

## ### Input:

[TLE] The Wikipedia page is about NBA conference finals. The Wikipedia section is about eastern conference finals. The table headers are: | year | champion | ... You need to populate the column: year. [SEED] The seed entity is <1971\_NBA\_playoffs>.

## ### Question:

The entity candidates are: <2003\_NBA\_playoffs>, <1982-83\_Washington\_Bullets\_season>, <2004\_NBA\_playoffs>, <Philadelphia\_76ers>, <1983-84\_Washington\_Bullets\_season>, <1952\_NBA\_playoffs>, ...

### Response: <1972\_NBA\_playoffs>, <1973\_NBA\_playoffs>, <1974\_NBA\_playoffs>, <1975\_NBA\_playoffs>, <1976\_NBA\_playoffs>, ...

## ### Instruction:

This is a **hierarchical table question answering** task. The goal for this task is to answer the given question based on the given table. The table might be hierarchical.

## ### Input:

[TLE] The table caption is department of defense obligations for research, development, test, and evaluation, by agency: 2015-18. [TAB] | agency | 2015 | 2016 | ... [SEP] | department of defense | department of defense | ... [SEP] | rdt&e | 61513.5 | ... [SEP] | total research | 6691.5 | ... [SEP] | basic research | 2133.4 | ... [SEP] | defense advanced research projects agency | ...

## ### Question:

How many dollars are the difference for basic research of defense advanced research projects agency increase between 2016 and 2018?

### Response: 80.3.

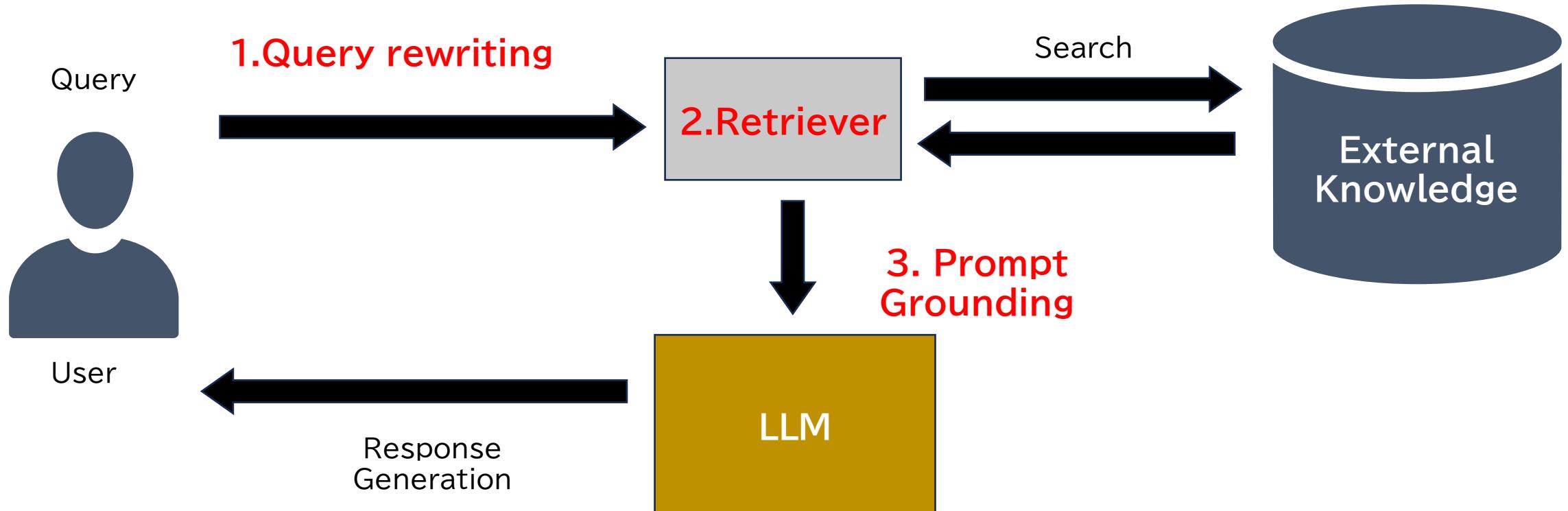
# Q&A

# On the Use of Large Language Models for Table Tasks

- Retrieval-augmented generation (RAG)

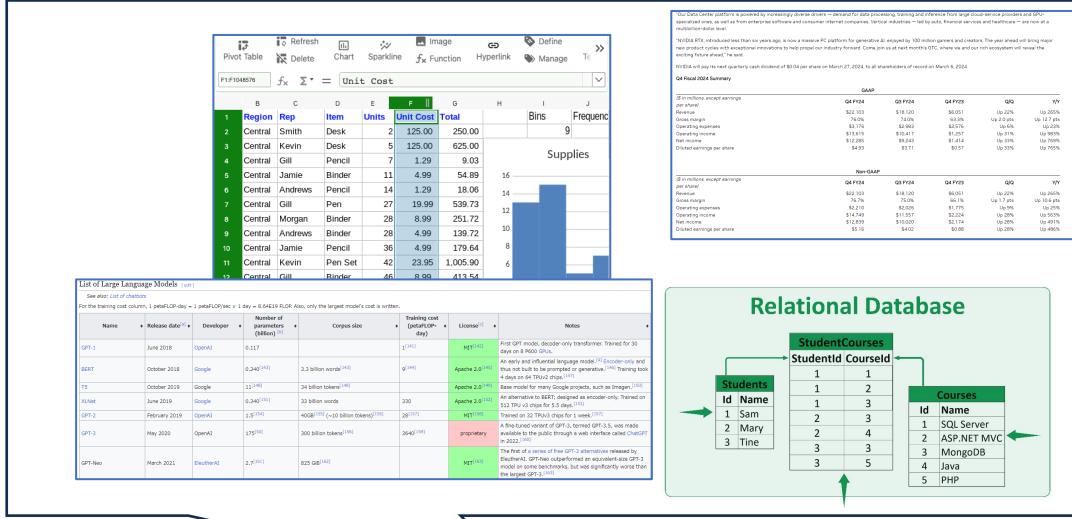
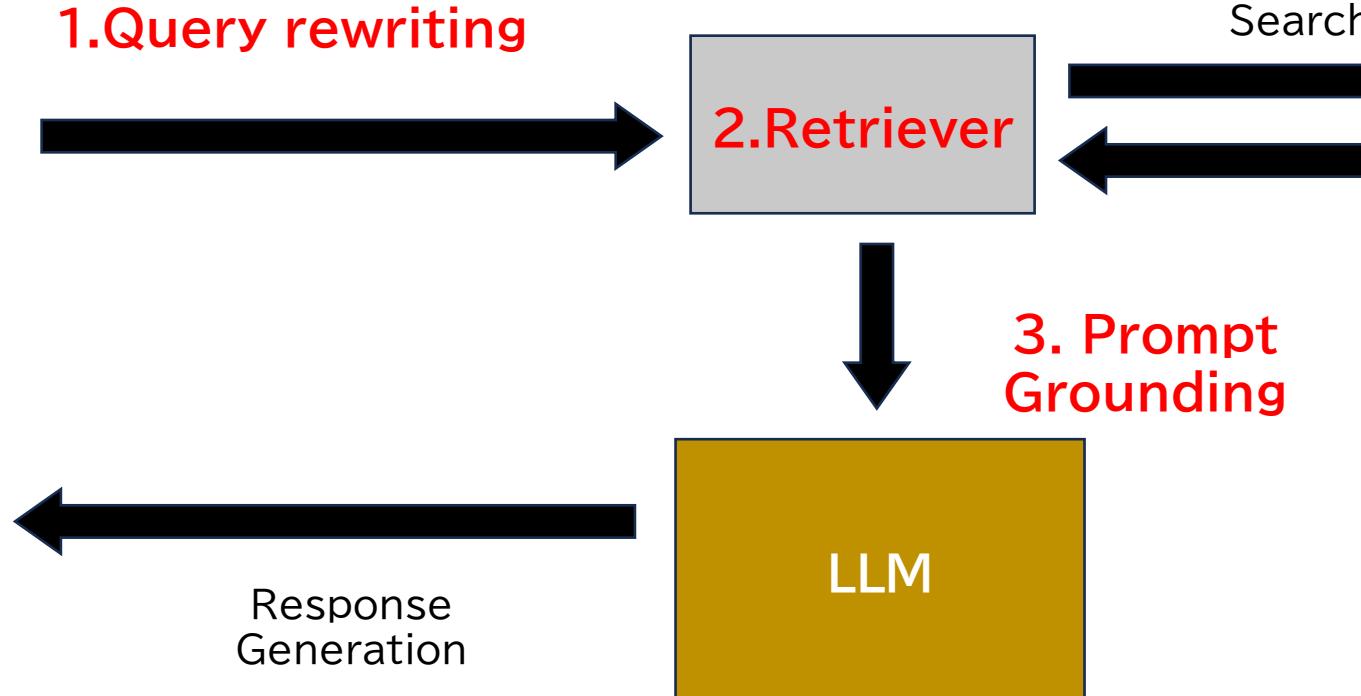
# What is RAG?

- LLM can not be trained frequently
- RAG is an efficient way to on-demand get external knowledge
  - Up-to-date knowledge
  - Less hallucination with retrieved context



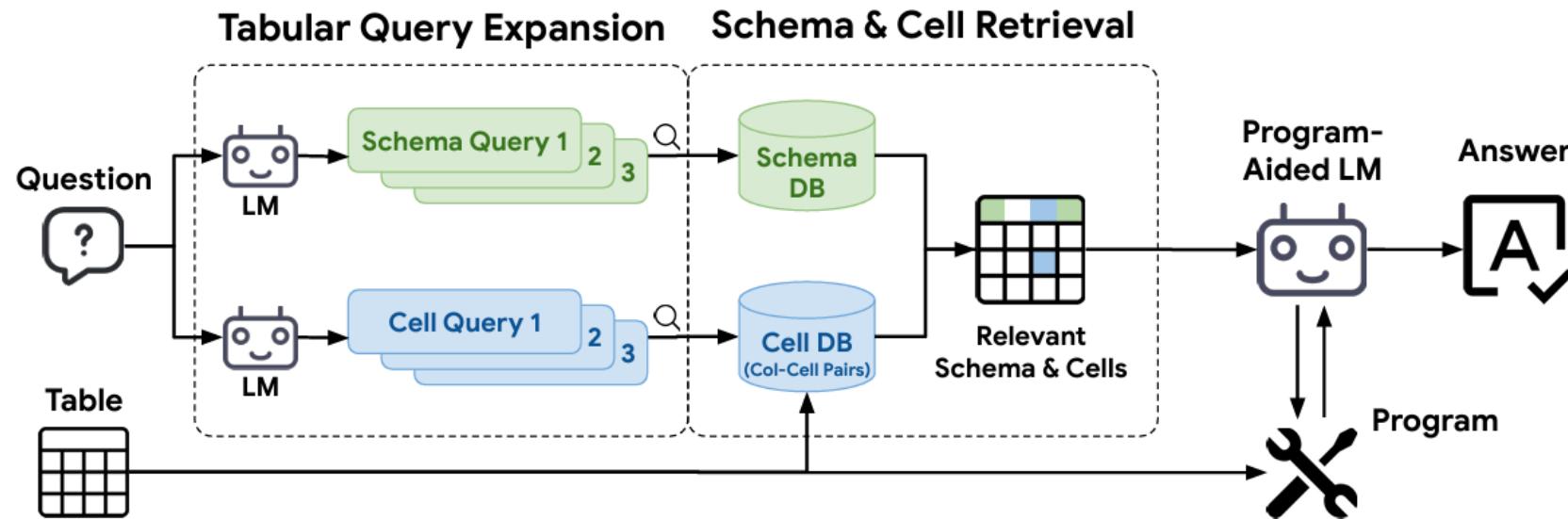
# Motivation of RAG with table

- Knowledge is also in tables!



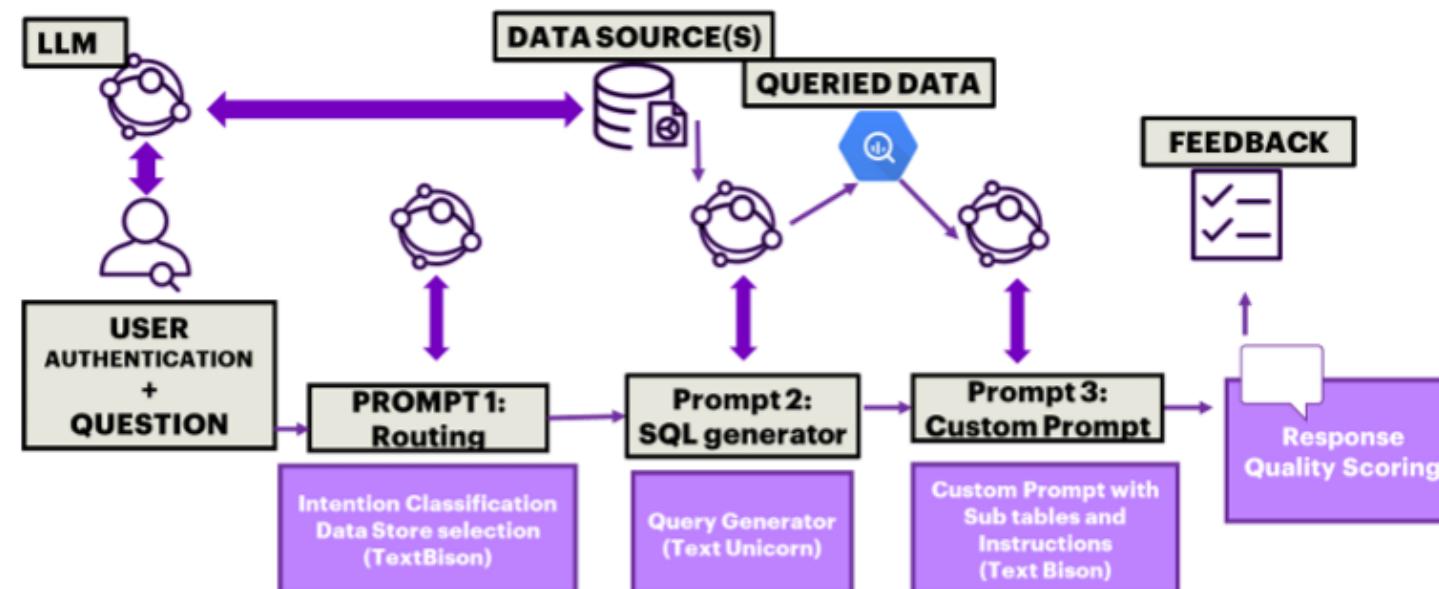
# Query rewriting & Prompt Grounding

- TableRAG [arxiv24]: for table QA
  - Query rewriting
    - NL query to Schema Query and Cell Query
    - NL query to SQL query
  - Prompt Grounding
    - Python generation with <query, retrieved data>
      - Execute with python shell
      - Observe result then run another round of generation until getting final answer



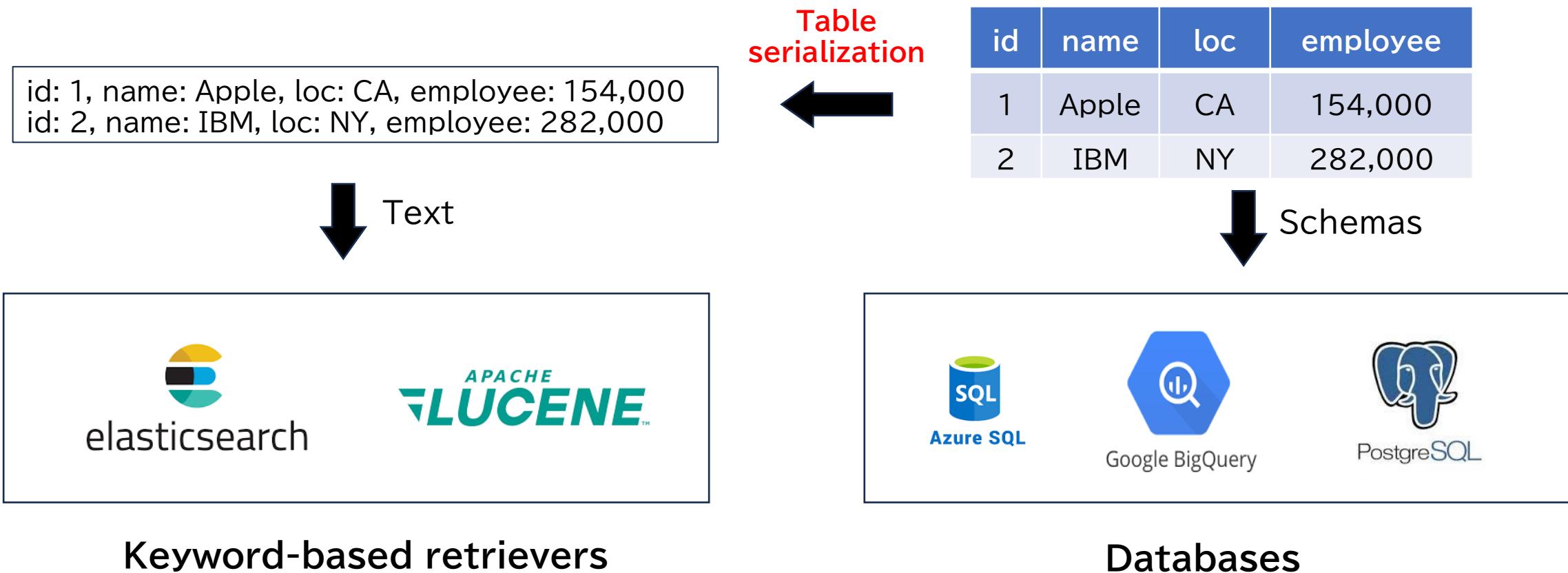
# Query rewriting & Prompt Grounding

- ERATTA [BigData24]: for table QA
  - Query rewriting
    - Routing: select target table
    - SQL generation: query to target table
  - Grounding
    - SQL generation with  $\langle$ query, retrieved data $\rangle$ 
      - Execute with BigQuery
      - Result feedback with 5 checks (entity, number, query, Regurgitation, Modifier)



# Table Retriever

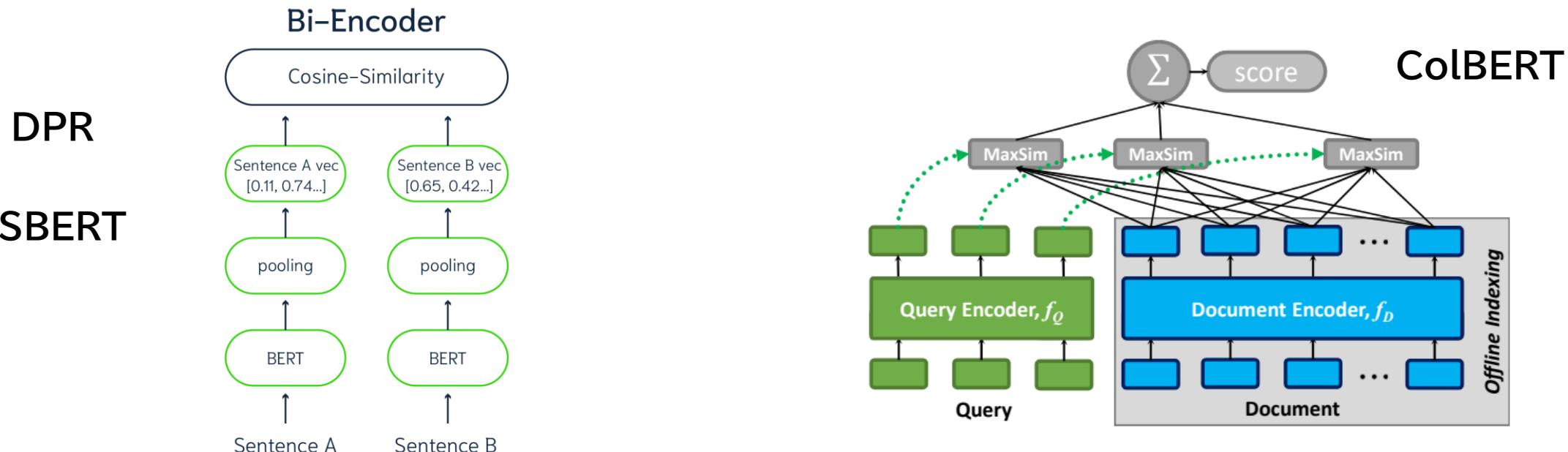
- Traditional retriever
  - Keyword-based (BM25) for text-based retriever
  - SQL for relational database.
    - ERATTA, TableRAG, DB-GPT[arxiv23] <https://arxiv.org/abs/2312.17449>



# Table Retriever

- Dense Retriever: LM encoder + vector database
  - Textual model: serialize table to text -> train text LM encoder

|                   | Model-architecture | serialization     |
|-------------------|--------------------|-------------------|
| T-RAG [arxiv22]   | DPR                | rows              |
| Deep-join[VLDB23] | Sentence-BERT      | columns           |
| LI-RAGE [ACL23]   | ColBERT            | table             |
| ITR [ACL23]       | DPR                | located row & col |



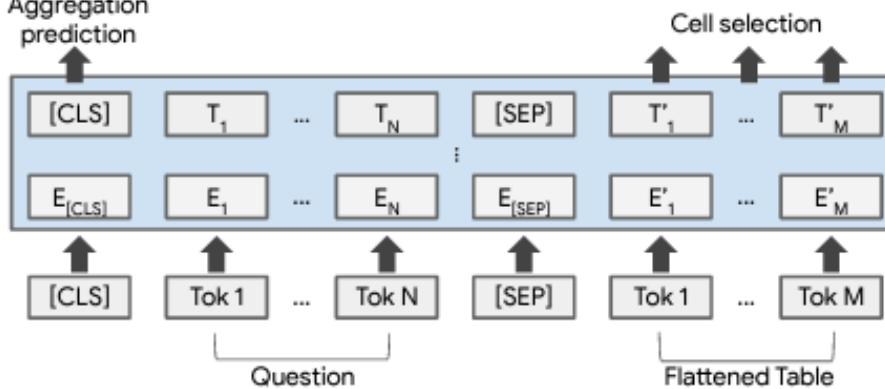
# Table Retriever

- Dense Retriever: LM encoder + vector database
  - Table-specific model: directly train the table with table-specific LM encoder
    - query -> text encoder
    - table -> table-specific encoder

| op    | $P_s(op)$ | compute(op, $P_s$ , T)       |
|-------|-----------|------------------------------|
| NONE  | 0         | -                            |
| COUNT | 0.1       | .9 + .9 + .2 = 2             |
| SUM   | 0.8       | .9x37 + .9x31 + .2x15 = 64.2 |
| AVG   | 0.1       | $64.2 \div 2 = 32.1$         |

$S_{pred} = .1 \times 2 + .8 \times 64.2 + .1 \times 32.1 = 54.8$

Aggregation prediction



TAPAS

| Rank | ... | Days | $P_s$ |
|------|-----|------|-------|
| 1    | ... | 37   | 0.9   |
| 2    | ... | 31   | 0.9   |
| 3    | ... | 17   | 0     |
| 4    | ... | 15   | 0.2   |
| ...  | ... | ...  | 0     |

In which city did Piotr's last 1st place finish occur?

|       | Year | Venue   | Position | Event                     |
|-------|------|---------|----------|---------------------------|
| $R_1$ | 2003 | Tampere | 3rd      | EU Junior Championship    |
| $R_2$ | 2005 | Erfurt  | 1st      | EU U23 Championship       |
| $R_3$ | 2005 | Izmir   | 1st      | Universiade               |
| $R_4$ | 2006 | Moscow  | 2nd      | World Indoor Championship |
| $R_5$ | 2007 | Bangkok | 1st      | Universiade               |

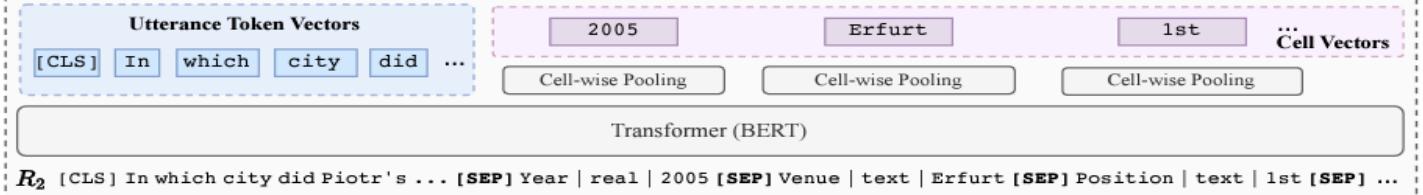
Selected Rows as Content Snapshot :  $\{R_2, R_3, R_5\}$

(A) Content Snapshot from Input Table



(C) Vertical Self-Attention over Aligned Row Encodings

(B) Per-row Encoding (for each row in content snapshot, using  $R_2$  as an example)



$R_2$  [CLS] In which city did Piotr's ... [SEP] Year | real | 2005 [SEP] Venue | text | Erfurt [SEP] Position | text | 1st [SEP] ...

TaBERT

# Table Retriever

- Do we need a table-specific model for table retriever?
- Table Retrieval May Not Necessitate Table-specific Model Design [ACL workshop22]
  - Table-specific(DTR) do not have significant improvement than text retrievers (DPR)
- Open-WikiTable [ACL23]
  - Same observation

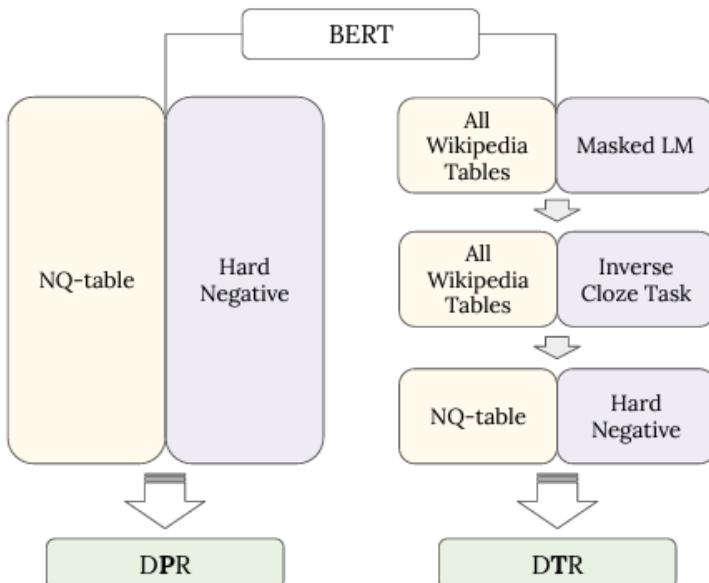


Figure 3: Comparison of DPR and DTR training.

| Model        | Retrieval Accuracy |              |              |              |              |
|--------------|--------------------|--------------|--------------|--------------|--------------|
|              | @1                 | @5           | @10          | @20          | @50          |
| DTR (medium) | 62.32              | 82.51        | 86.75        | 91.51        | 94.26        |
| DTR (large)  | 63.98              | 84.27        | <b>89.65</b> | <b>93.48</b> | <b>95.65</b> |
| BERT-table   | 60.97              | 79.81        | 85.51        | 88.20        | 91.62        |
| DPR          | 57.04              | 80.54        | 86.13        | 89.54        | 92.34        |
| DPR-table    | <b>67.91</b>       | <b>84.89</b> | 88.72        | 90.58        | 92.86        |

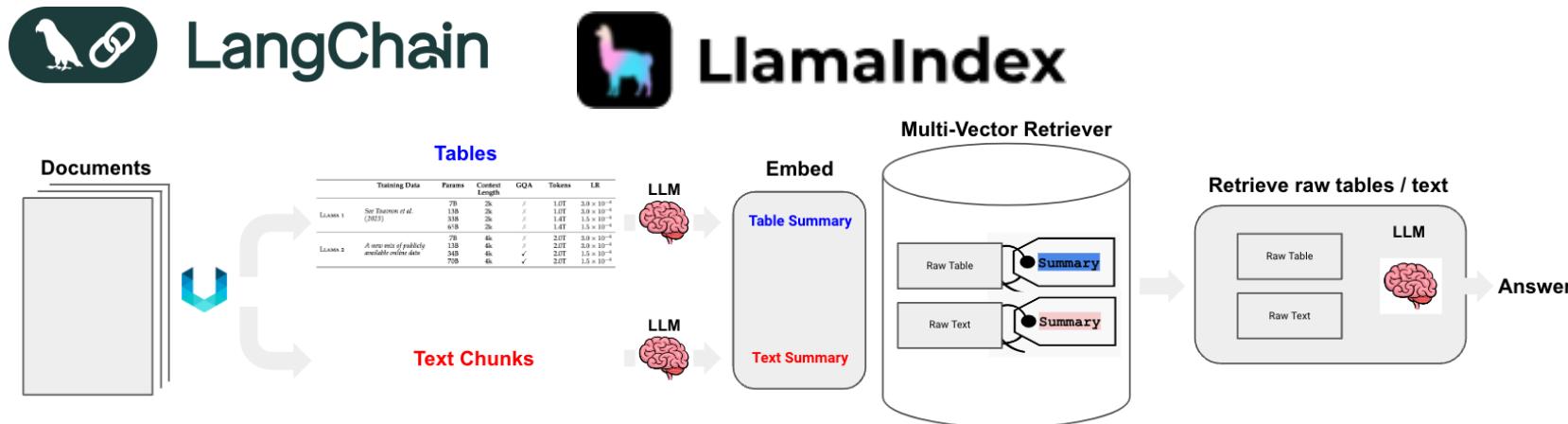
NQ-tables

| Encoder | Data  | k=5 k=10 k=20    |       |      |      |      |
|---------|-------|------------------|-------|------|------|------|
|         |       | Text             | Table | k=5  | k=10 | k=20 |
| BERT    | BERT  | Original         |       | 6.6  | 8.0  | 10.3 |
|         |       | Decontextualized |       | 45.5 | 52.9 | 59.7 |
|         |       | Paraphrased      |       | 42.2 | 48.9 | 56.1 |
| TAPAS   | TAPAS | Original         |       | 25.0 | 34.1 | 45.1 |
|         |       | Decontextualized |       | 91.6 | 96.0 | 97.8 |
|         |       | Paraphrased      |       | 89.5 | 95.0 | 97.3 |

Open-WikitaTables

# Table Retriever (industry-view)

- Entry: Keyword + traditional retriever
  - Elasticsearch: <https://www.elastic.co/search-labs/blog/rag-with-llamaIndex-and-elasticsearch>
  - RAG with Azure SQL: <https://devblogs.microsoft.com/azure-sql/rag-with-azure-sql/>
  - LlamaIndex: <https://docs.llamaindex.ai/en/v0.10.18/api reference/indices/table.html>
- Advanced: LLM parse + textual embedding + vector database



<https://blog.langchain.dev/semi-structured-multi-modal-rag/>  
<https://blog.langchain.dev/benchmarking-rag-on-tables/>  
[https://docs.llamaindex.ai/en/stable/llama\\_cloud/llama\\_parse/](https://docs.llamaindex.ai/en/stable/llama_cloud/llama_parse/)



102

- Deluxe (not recommended) : Human parse + textual embedding + vector database

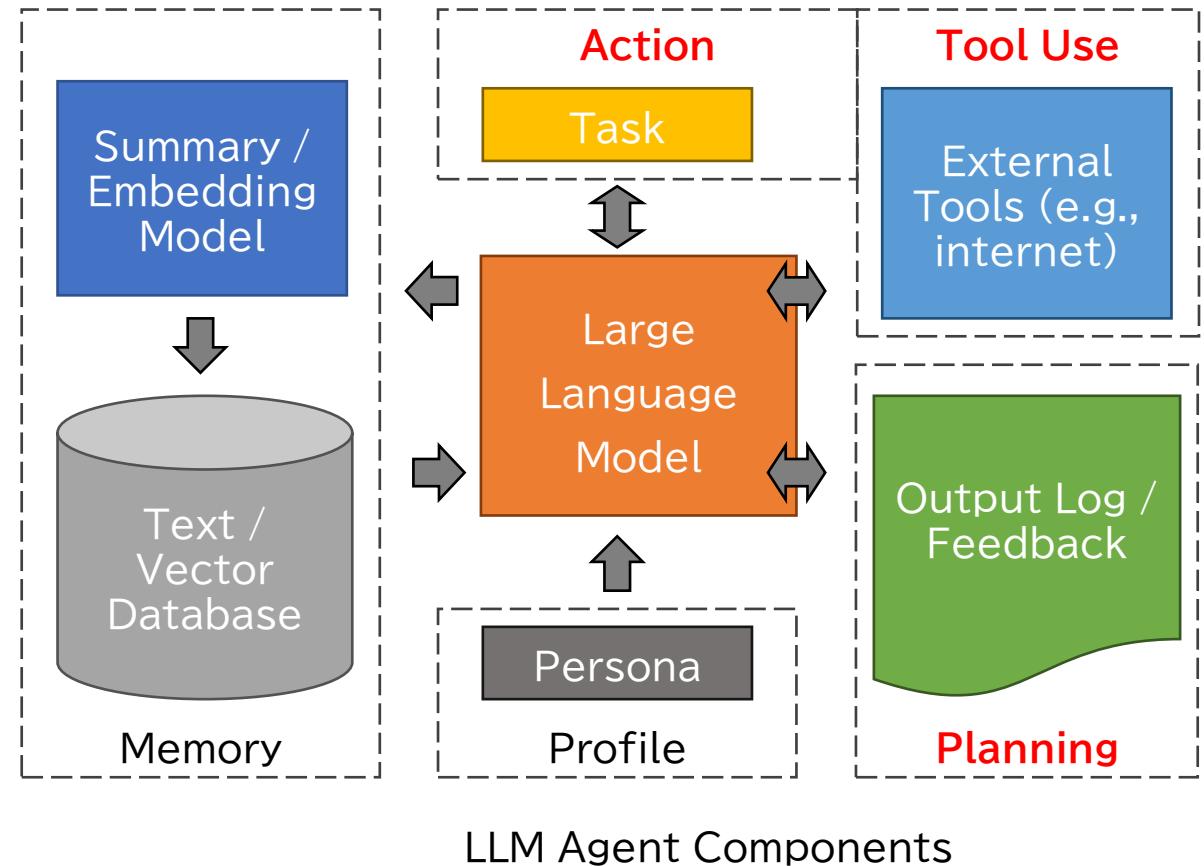
# Q&A

# On the Use of Large Language Models for Table Tasks

- LLM agents

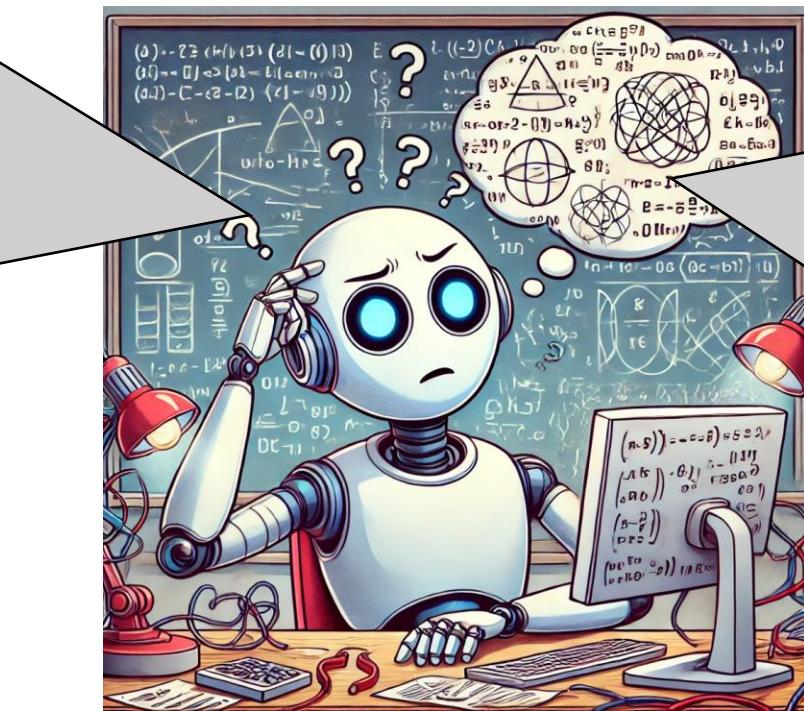
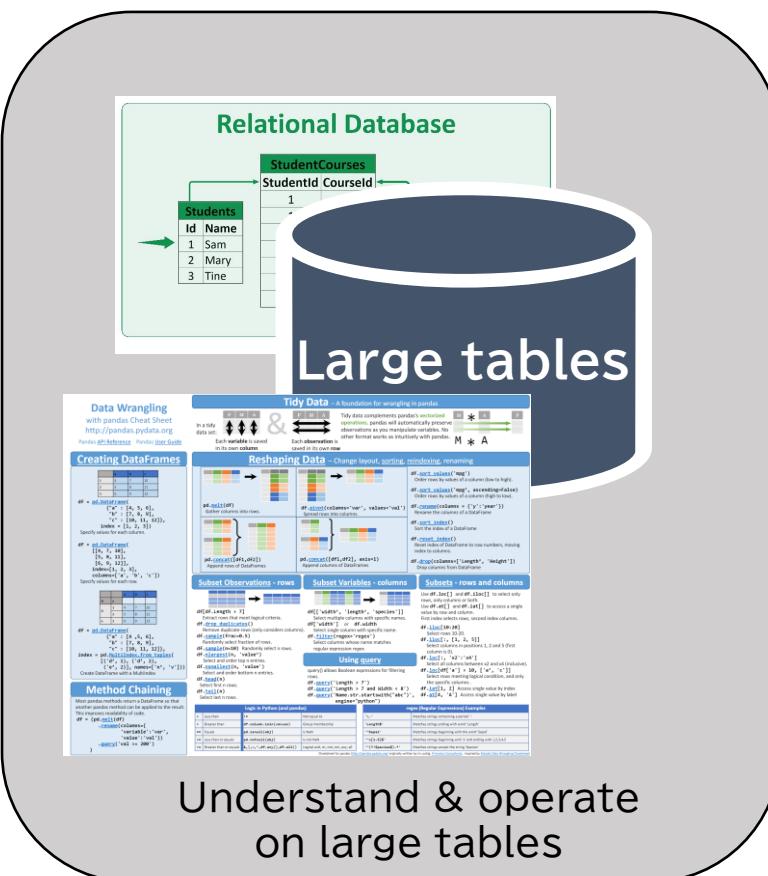
# What is LLM agent

- **Planning**
  - Task decomposition
  - Task selection
  - Reflection & refinement
- **Action & Tool use**
  - Design an action space
  - Interactive with external tools
- **Profile**
  - Role play, prompt.
  - “You are an expert of table/SQL/Pandas...”
- **Memory**
  - Short-term memory (log in a dialog)
  - Long-term memory (RAG with outside knowledge)



# Motivation of agent for table

- One-time inference can not solve hard problems, need the “divide and concur”
- Table is complicate, only a part of contents is useful for a specific question.
- Answer a question need multi-hop reasoning.
- Need external tools
  - Numerical OLAP, aggregation, SQL, python, Pandas ...



|                     |                     |
|---------------------|---------------------|
| square beads        | \$2.97 per kilogram |
| oval beads          | \$3.41 per kilogram |
| flower-shaped beads | \$2.18 per kilogram |
| star-shaped beads   | \$1.95 per kilogram |
| heart-shaped beads  | \$1.52 per kilogram |
| spherical beads     | \$3.42 per kilogram |
| rectangular beads   | \$1.97 per kilogram |

| Sandwich sales      |      |           |
|---------------------|------|-----------|
| Shop                | Tuna | Egg salad |
| City Cafe           | 6    | 5         |
| Sandwich City       | 3    | 12        |
| Express Sandwiches  | 7    | 17        |
| Sam's Sandwich Shop | 1    | 6         |
| Kelly's Subs        | 3    | 4         |

Mathematical reasoning from tables

**Question:** If Tracy buys 5 kilograms of spherical beads, 4 kilograms of star-shaped beads, and 3 kilograms of flower-shaped beads, how much will she spend? (unit: \$)

**Answer:** 31.44

**Solution:**  
Find the cost of the spherical beads. Multiply:  $\$3.42 \times 5 = \$17.10$ .  
Find the cost of the star-shaped beads. Multiply:  $\$1.95 \times 4 = \$7.80$ .  
Find the cost of the flower-shaped beads. Multiply:  $\$2.18 \times 3 = \$6.54$ .  
Now find the total cost by adding:  $\$17.10 + \$7.80 + \$6.54 = \$31.44$ .  
She will spend \$31.44.

**Question:** As part of a project for health class, Cara surveyed local delis about the kinds of sandwiches sold. Which shop sold fewer sandwiches, Sandwich City or Express Sandwiches?

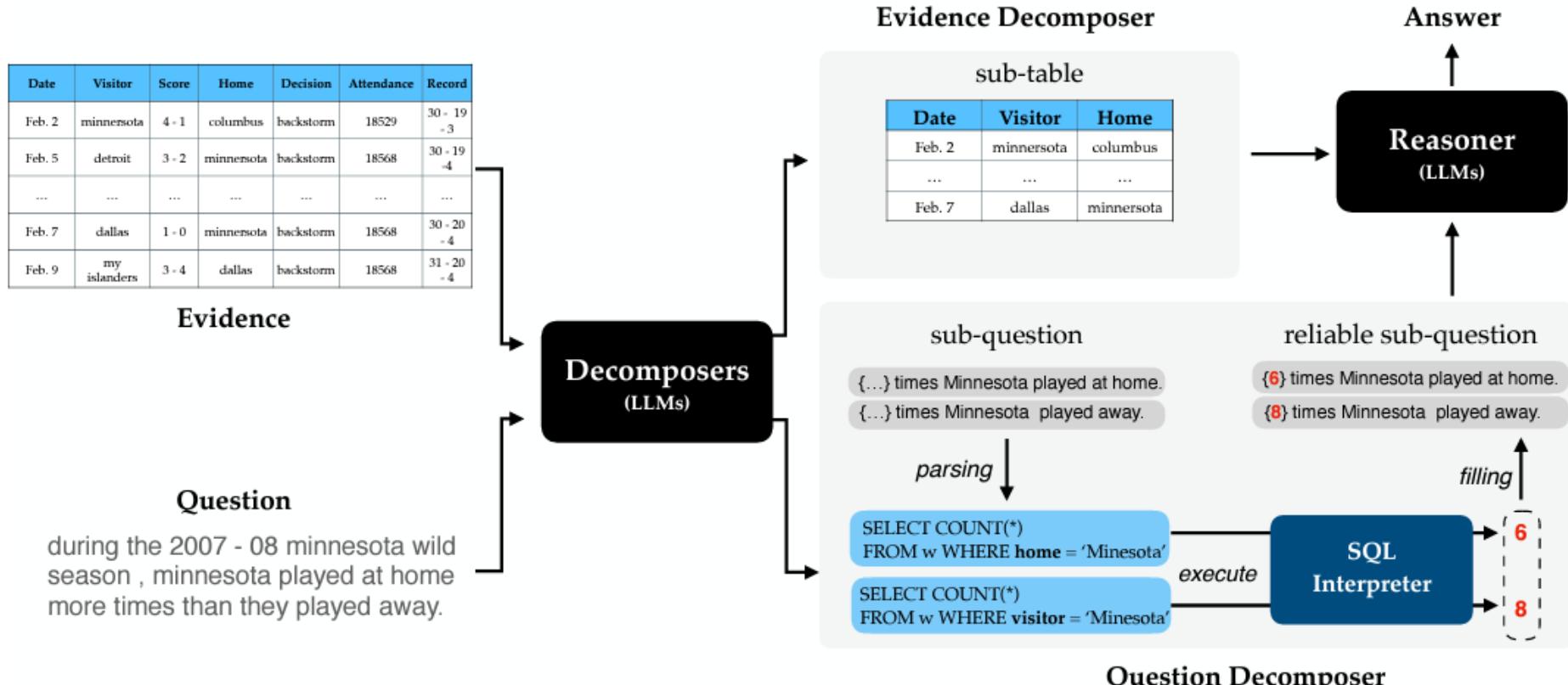
**Options:** (A) Sandwich City (B) Express Sandwiches

**Answer:** (A) Sandwich City

**Solution:**  
Add the numbers in the Sandwich City row. Then, add the numbers in the Express Sandwiches row.  
Sandwich City:  $3 + 12 = 15$ . Express Sandwiches:  $7 + 17 = 24$ .  
 $15$  is less than  $24$ . Sandwich City sold fewer sandwiches.

# Planning:

- Task decomposition
  - Datar [SIGIR23]: Table QA
  - Question decomposition and table decomposition
    - sub question -> SQL -> query sub-table -> answer



# Planning:

- Task decomposition & Reflection

- DIN-SQL [NeurIPS23]: Text-to-SQL
- Predefine the decomposition into  $1 + 3 = 4$  sub tasks
  - Schema linking -> Generate 3-levels SQL with 3 sub tasks
- Self-correction the final answer

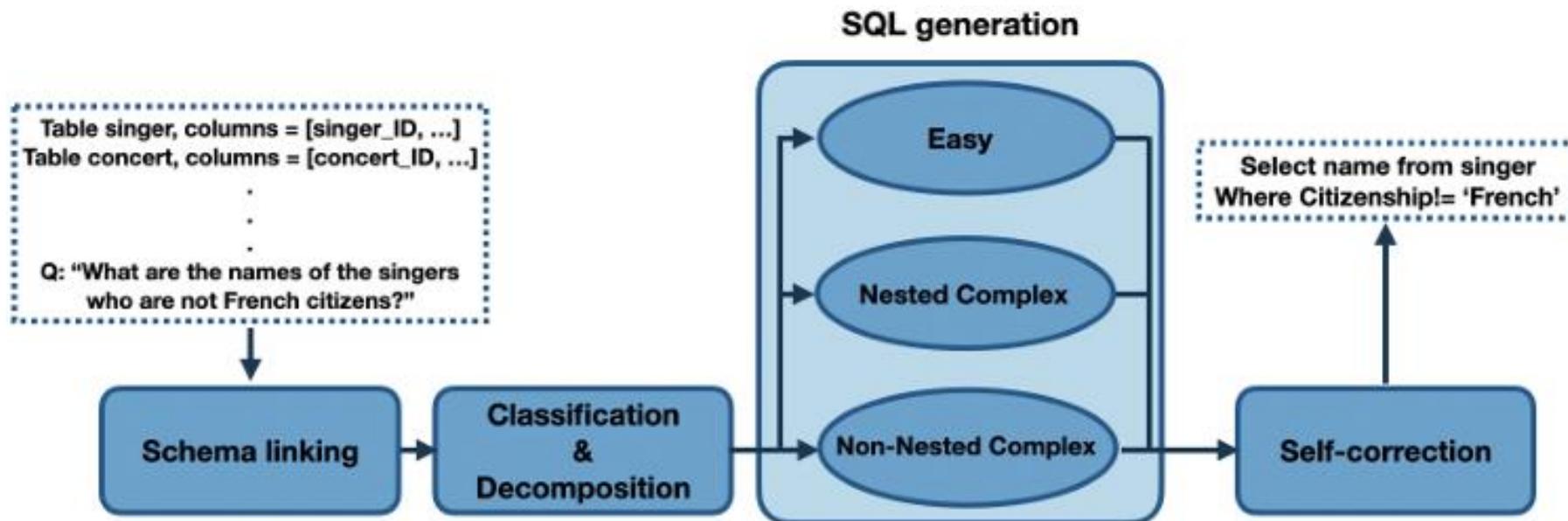
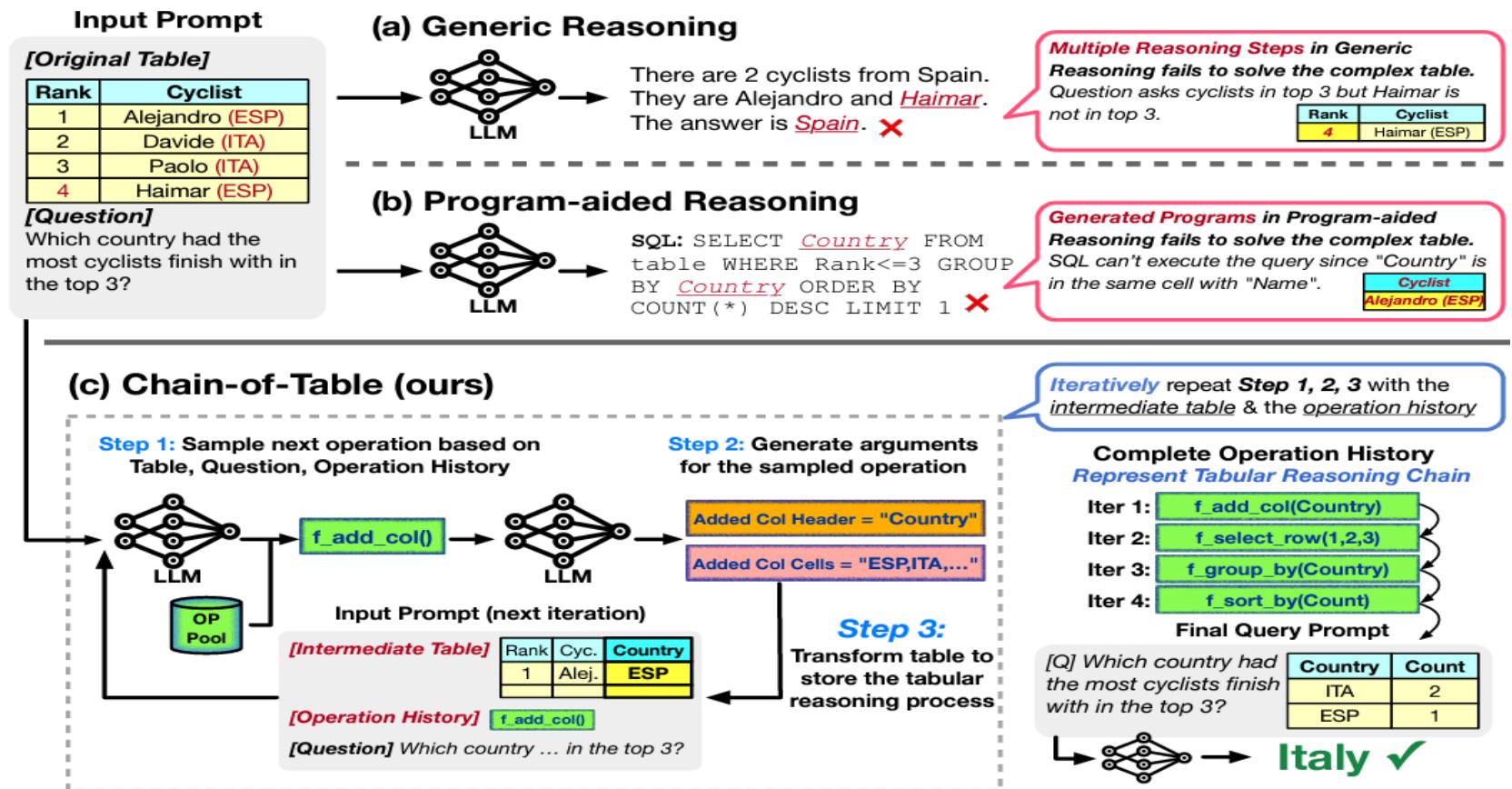


Figure 2: An overview of the proposed methodology including all four modules

# Planning:

- Task decomposition & Task selection & Reflection
  - Chain-of-table [ICLR24]: TableQA, TabFact
  - Task decomposition: select action and generate arguments
  - Plan selection: DynamicPlan to select next action based on the current state



# Action & Tool use:

- Action Space Design

|                | Task                 | Action Space   |
|----------------|----------------------|--|
| Chain-of-Table | TableQA              | Table actions:<br>[add_col, select_row, select_col, group_by, sort_by]   |
| ToolQA         | TableQA              | Tools:<br>[text, database, math, graph, code, system]  |
| SheetCopilot   | Table Transformation | Virtual APIS of 40+ atomic actions<br>Categories: [manipulation, management, formatting, Charts, Pivot Table, Formula] |

# Action & Tool use:

- External tool use

|                | Task                 | SQL | Python | Retriever | Database | Math | Graph |
|----------------|----------------------|-----|--------|-----------|----------|------|-------|
| PyAgent        | TableQA              |     | ○      |           |          |      |       |
| ToolQA         | TableQA              | ○   | ○      | ○         | ○        | ○    | ○     |
| ReAcTable      | TableQA              | ○   | ○      |           |          |      |       |
| LEVER          | TableQA              | ○   |        |           |          |      |       |
| Binder         | TableQA, TableFact   | ○   | ○      | ○         |          |      |       |
| Chain-of-Table | TableQA, TableFact   |     | ○      |           |          |      |       |
| ToolWriter     | TableQA              | ○   | ○      |           |          |      |       |
| SheetAgent     | Table Transformation | ○   | ○      | ○         |          |      |       |
| AutoTQA        | TableQA, TableFact   | ○   |        | ○         | ○        |      |       |
| TableRAG       | TableQA, TableFact   |     | ○      | ○         |          |      |       |

# Action & Tool use:

- SheetCopilot: planning, action & tool use (generate API call) for table transformation.

Instrucion: Please highlight Sales between 200 and 500.

**Target Spreadsheet**

Context: My workbook records many invoices made on different dates.

| No. | Sales Rep | Product | Price    | Units   | Sales |          |
|-----|-----------|---------|----------|---------|-------|----------|
| 1   | 10500     | Joe     | Majestic | \$30.00 | 25    | \$750.00 |
| 2   | 10501     | Moe     | Majestic | \$30.00 | 9     | \$270.00 |
| 3   | 10501     | Moe     | Quad     | \$32.00 | 21    | \$672.00 |
| 4   | 10501     | Moe     | Alpine   | \$22.00 | 7     | \$154.00 |
| 5   | 10501     | Moe     | Carlota  | \$25.00 | 11    | \$275.00 |
| 6   | 10502     | Moe     | Majestic | \$30.00 | 5     | \$150.00 |
| 7   | 10502     | Moe     | Carlota  | \$25.00 | 25    | \$625.00 |
| 8   | 10503     | Chin    | Carlota  | \$25.00 | 21    | \$525.00 |
| 9   | 10503     | Chin    | Alpine   | \$22.00 | 16    | \$352.00 |

## Planning

Generate Prompt  
`New_prompt = Concatenate(Base_prompt,  
Detailed_doc)`

Get Detailed\_Doc  
Filter:  
...

Revised Plan  
Step 1. Filter the range by the criteria ">=200" and "<=500".  
Action API: `@Filter(source="A!A1:F10",  
fieldIndex=6, criteria=">=200,<=500")@`

## Final Action

`Filter(source="A!A1:F10", fieldIndex=6,  
criteria=">=200,<=500")`

## Full Plan and Execution

Step 1. Filter the range by the criteria ">=200" and "<=500".

Step 2. Set the fill color of the filtered cells to green.

Step 3. Remove the filter.

Action API:  
`Filter(source="A!A1:F10",  
fieldIndex=6,  
criteria=">=200,<=500")`

Action API:  
`SetFormat(source="A!A1:F10",  
fillColor="pale_blue")`

Action API:  
`DeleteFilter()`

| No. | Sales Rep | Product | Price    | Units   | Sales |          |
|-----|-----------|---------|----------|---------|-------|----------|
| 1   | 10500     | Joe     | Majestic | \$30.00 | 25    | \$750.00 |
| 2   | 10501     | Moe     | Majestic | \$30.00 | 9     | \$270.00 |
| 3   | 10501     | Moe     | Quad     | \$32.00 | 21    | \$672.00 |
| 4   | 10501     | Moe     | Alpine   | \$22.00 | 7     | \$154.00 |
| 5   | 10501     | Moe     | Carlota  | \$25.00 | 11    | \$275.00 |
| 6   | 10502     | Moe     | Majestic | \$30.00 | 5     | \$150.00 |
| 7   | 10502     | Moe     | Carlota  | \$25.00 | 25    | \$625.00 |
| 8   | 10503     | Chin    | Carlota  | \$25.00 | 21    | \$525.00 |
| 9   | 10503     | Chin    | Alpine   | \$22.00 | 16    | \$352.00 |

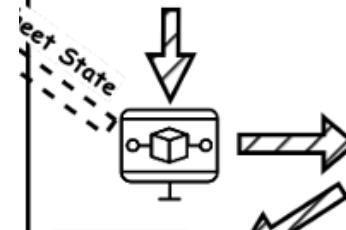
| No. | Sales Rep | Product | Price    | Units   | Sales |          |
|-----|-----------|---------|----------|---------|-------|----------|
| 1   | 10500     | Joe     | Majestic | \$30.00 | 25    | \$750.00 |
| 2   | 10501     | Moe     | Majestic | \$30.00 | 9     | \$270.00 |
| 3   | 10501     | Moe     | Quad     | \$32.00 | 21    | \$672.00 |
| 4   | 10501     | Moe     | Alpine   | \$22.00 | 7     | \$154.00 |
| 5   | 10501     | Moe     | Carlota  | \$25.00 | 11    | \$275.00 |
| 6   | 10502     | Moe     | Majestic | \$30.00 | 5     | \$150.00 |
| 7   | 10502     | Moe     | Carlota  | \$25.00 | 25    | \$625.00 |
| 8   | 10503     | Chin    | Carlota  | \$25.00 | 21    | \$525.00 |
| 9   | 10503     | Chin    | Alpine   | \$22.00 | 16    | \$352.00 |

| No. | Sales Rep | Product | Price    | Units   | Sales |          |
|-----|-----------|---------|----------|---------|-------|----------|
| 1   | 10500     | Joe     | Majestic | \$30.00 | 25    | \$750.00 |
| 2   | 10501     | Moe     | Majestic | \$30.00 | 9     | \$270.00 |
| 3   | 10501     | Moe     | Quad     | \$32.00 | 21    | \$672.00 |
| 4   | 10501     | Moe     | Alpine   | \$22.00 | 7     | \$154.00 |
| 5   | 10501     | Moe     | Carlota  | \$25.00 | 11    | \$275.00 |
| 6   | 10502     | Moe     | Majestic | \$30.00 | 5     | \$150.00 |
| 7   | 10502     | Moe     | Carlota  | \$25.00 | 25    | \$625.00 |
| 8   | 10503     | Chin    | Carlota  | \$25.00 | 21    | \$525.00 |
| 9   | 10503     | Chin    | Alpine   | \$22.00 | 16    | \$352.00 |

## Tool use

## #3 Acting Stage

Execute generated actions in Simulation Environment



| No. | Sales Rep | Product | Price    | Units   | Sales |          |
|-----|-----------|---------|----------|---------|-------|----------|
| 1   | 10500     | Joe     | Majestic | \$30.00 | 25    | \$750.00 |
| 2   | 10501     | Moe     | Majestic | \$30.00 | 9     | \$270.00 |
| 3   | 10501     | Moe     | Quad     | \$32.00 | 21    | \$672.00 |
| 4   | 10501     | Moe     | Alpine   | \$22.00 | 7     | \$154.00 |
| 5   | 10501     | Moe     | Carlota  | \$25.00 | 11    | \$275.00 |
| 6   | 10502     | Moe     | Majestic | \$30.00 | 5     | \$150.00 |
| 7   | 10502     | Moe     | Carlota  | \$25.00 | 25    | \$625.00 |
| 8   | 10503     | Chin    | Carlota  | \$25.00 | 21    | \$525.00 |
| 9   | 10503     | Chin    | Alpine   | \$22.00 | 16    | \$352.00 |

## Action

Go to Proposing Stage for the next step

# Action & Tool use:

- ReActTable [VLDB24]
  - Apply ReAct [ICLR24], reason and action concept to table QA task
  - During action, use SQL and python shell

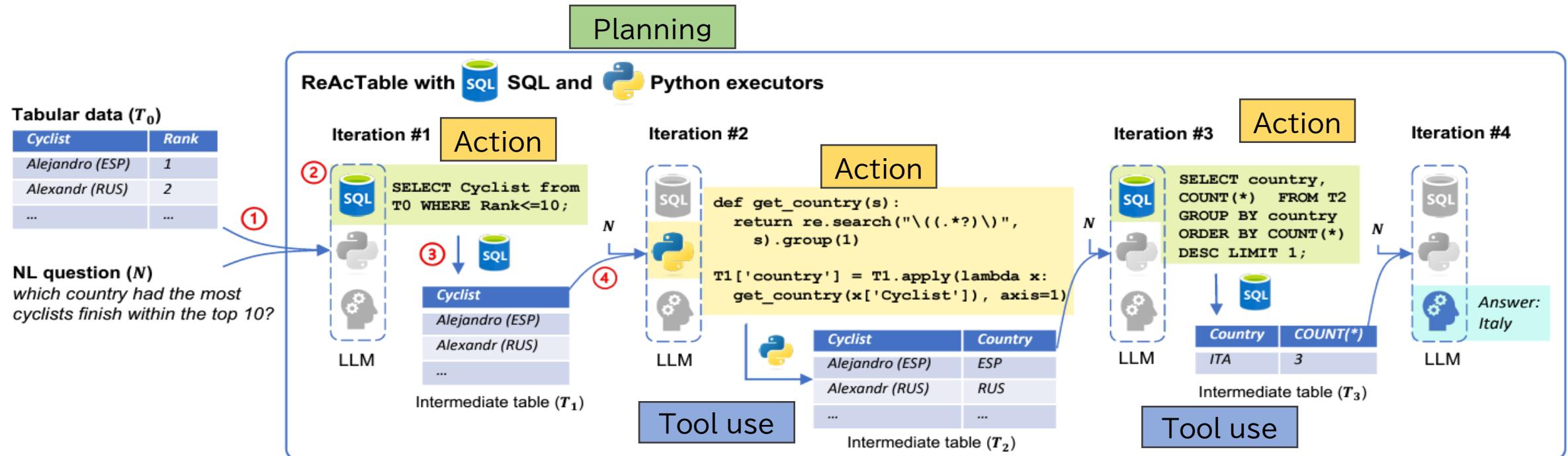
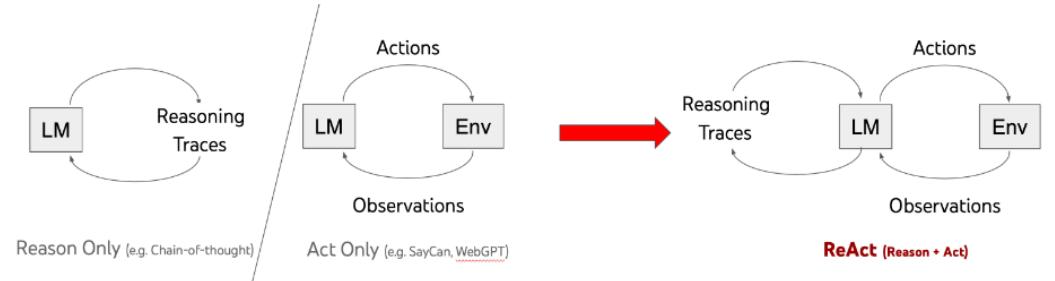
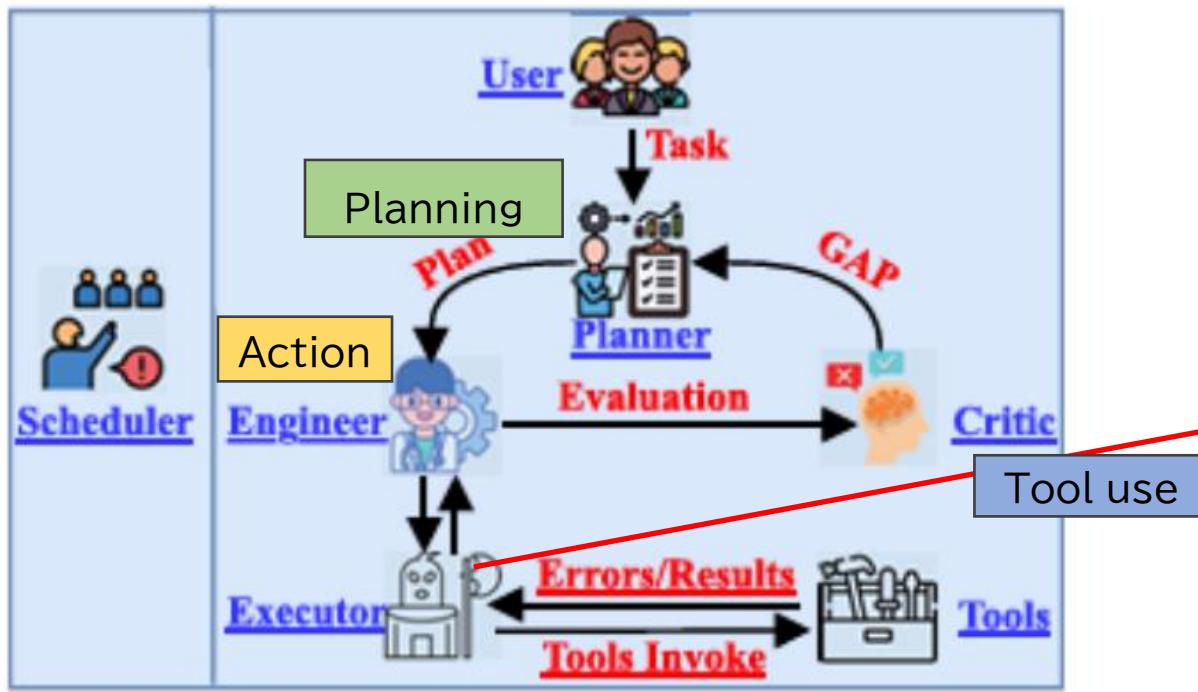


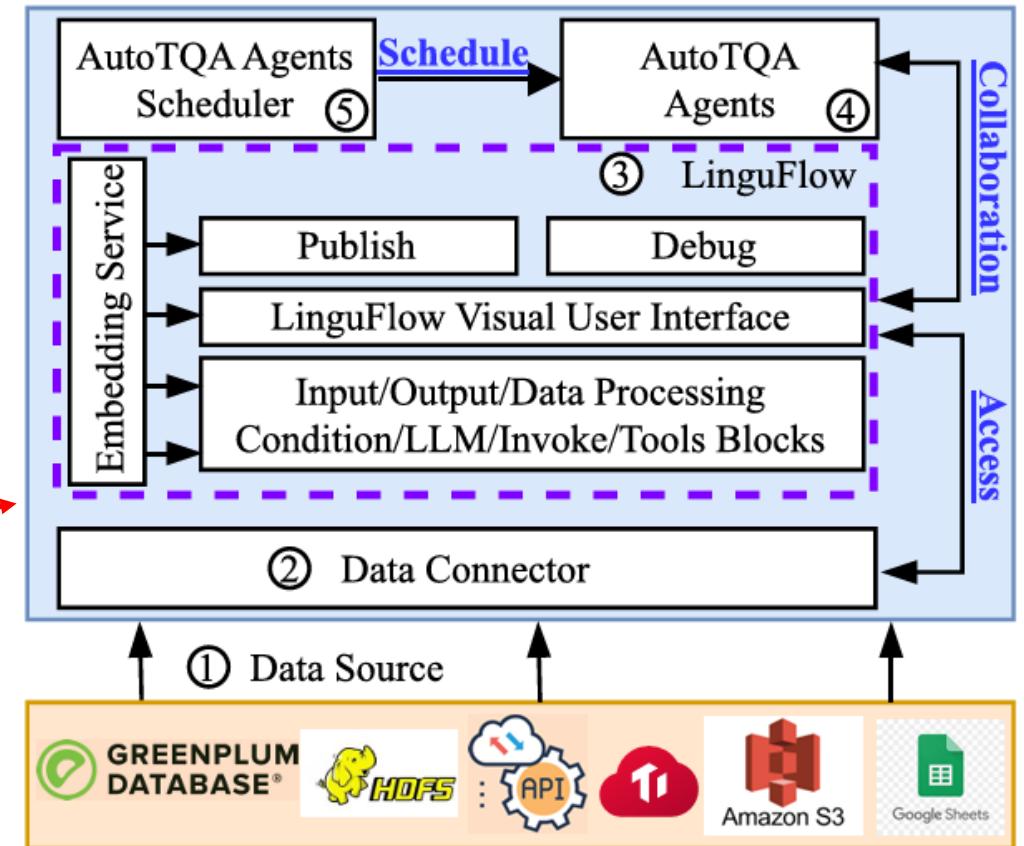
Figure 1: Overview of the ReAcTable framework with SQL and Python code executors.

# Action & Tool use:

- AutoTQA [VLDB24]: Use profile and separate action and tool use for agent cooperation.



Profiles (Planner, Engineer, Executor and Critic)



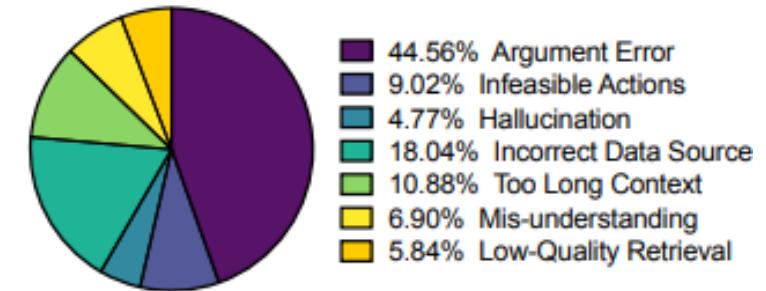
# Action & Tool use:

- ToolQA: a dataset to eval the ability of LLM use tool for table [arxiv23]
  - Include 13 tools in 8 domain

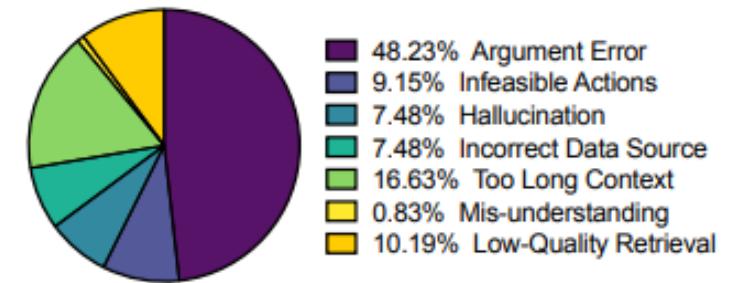
Table 2: Different tools in ToolQA.

| Tool Types     | # Tools | Tools   |
|----------------|---------|---|
| Text Tools     | 2       | Agenda Retriever, SciREX Retriever                          |
| Database Tools | 3       | Database Loader, Data Filter, Get Value                     |
| Math Tools     | 1       | WolframAlpha Calculator                                     |
| Graph Tools    | 4       | Graph Loader, Neighbour Checker, Node Checker, Edge Checker |
| Code Tools     | 2       | Python Interpreter, SQL Interpreter                         |
| System Tools   | 1       | Finish  |

| Error Type                   | Description  |
|------------------------------|--|
| Argument Error               | Incorrect parameters when calling tools                        |
| Incorrect Data Source        | Difficulty in identifying the proper reference corpora         |
| Innovation and Hallucination | Creative tool combinations leading to call a non-existent tool |
| Infeasible Actions           | Attempting to use non-existent tools                           |
| Too Long Context             | Exceeding model's length limitation                            |
| Misunderstanding             | Inability to comprehend tool execution feedback                |
| Low-Quality Retrieval        | Insufficient relevant information extracted from corpora       |



(a) Easy questions.



(b) Hard questions.

# Action & Tool use: Example of let LLM use program

You are working with a pandas dataframe in Python. The name of the dataframe is `df`. Your task is to use `python\_repl\_ast` to answer the question posed to you.

Profile

Tool description:

- `python\_repl\_ast`: A Python shell. Use this to execute python commands. Input should be a valid python command. When using this tool, ...

Tool use

Guidelines:

- **Aggregated Rows**: Be cautious of rows that aggregate data such as 'total', 'sum', or 'average'. Ensure these rows do not influence your results inappropriately.

- **Data Verification**: Before concluding the final answer, always verify that your observations align with the original table and question.

Strictly follow the given format to respond:

Planning

Question: the input question you must answer

Thought: you should always think about what to do to interact with `python\_repl\_ast`

Action: can **ONLY** be `python repl ast`

Action Input: the input code to the action

Observation: the result of the action

... (this Thought/Action/Action Input/Observation can repeat N times)

Thought: after verifying the table, observations, and the question, I am confident in the final answer

Final Answer: the final answer to the original input question (AnswerName1, AnswerName2...)

Notes for final answer:

- Ensure the final answer format is only "Final Answer: AnswerName1, AnswerName2..." form, no other form.
- Ensure the final answer is a number or entity names, as short as possible, without any explanation.
- Ensure to have a concluding thought that verifies the table, observations and the question before giving the final answer.

You are provided with a table regarding "[TITLE]". This is the result of `print(df.to\_markdown())`:

[TABLE]

**Note**: All cells in the table should be considered as `object` data type, regardless of their appearance.

Begin!

Question: [QUESTION]

"""

# Action & Tool use: Example of let LLM use APIs

Profile

You are now capable of calling a weather API to retrieve weather information.

The API format is:POST /weather{ "city": "{city\_name}"}

Tool use

When I ask for the weather in a city, respond by generating the API call in this format.

Let's see a few examples:

User: What's the weather like in Tokyo?

Response: API Call: POST /weather{ "city": "Tokyo"}

User: Tell me the weather in New York.

Response: API Call: POST /weather{ "city": "New York"}

User: Can you check the weather in London?

Response: API Call: POST /weather{ "city": "London"}

Now it's your turn to generate the correct API call when I ask about a city's weather.

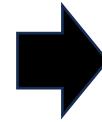
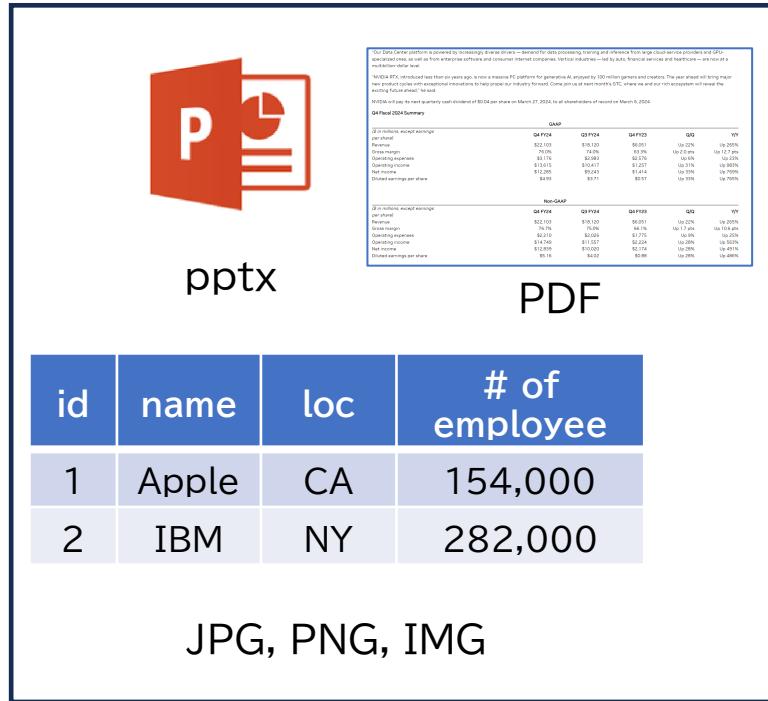
# Q&A

# On the Use of Large Language Models for Table Tasks

- Vision-language models (VLMs)

# What is VLM

- Input <prompt, image>, Output: response in text



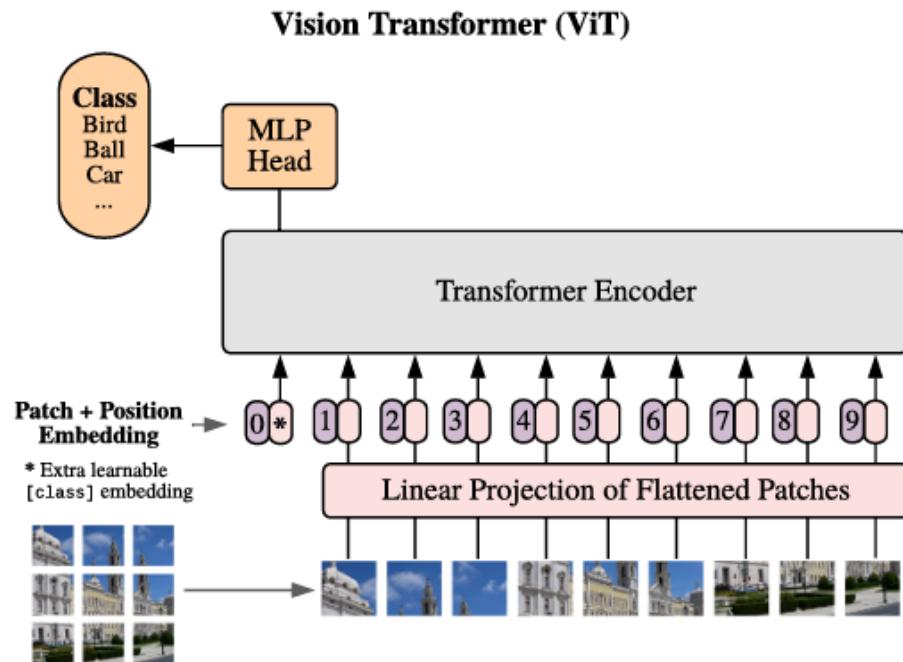
Response

The number of employees for Apple, as shown in the table, is 154,000.

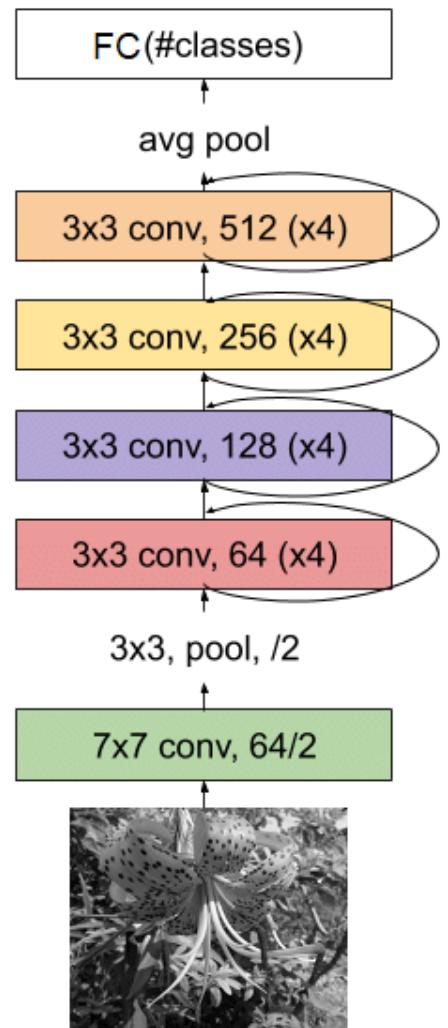
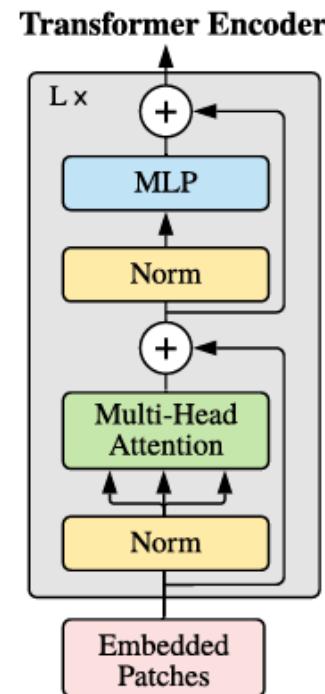
Prompt

what is the employee number of Apple?

# VIT (Vision Transformer) vs Resnet (CNN)

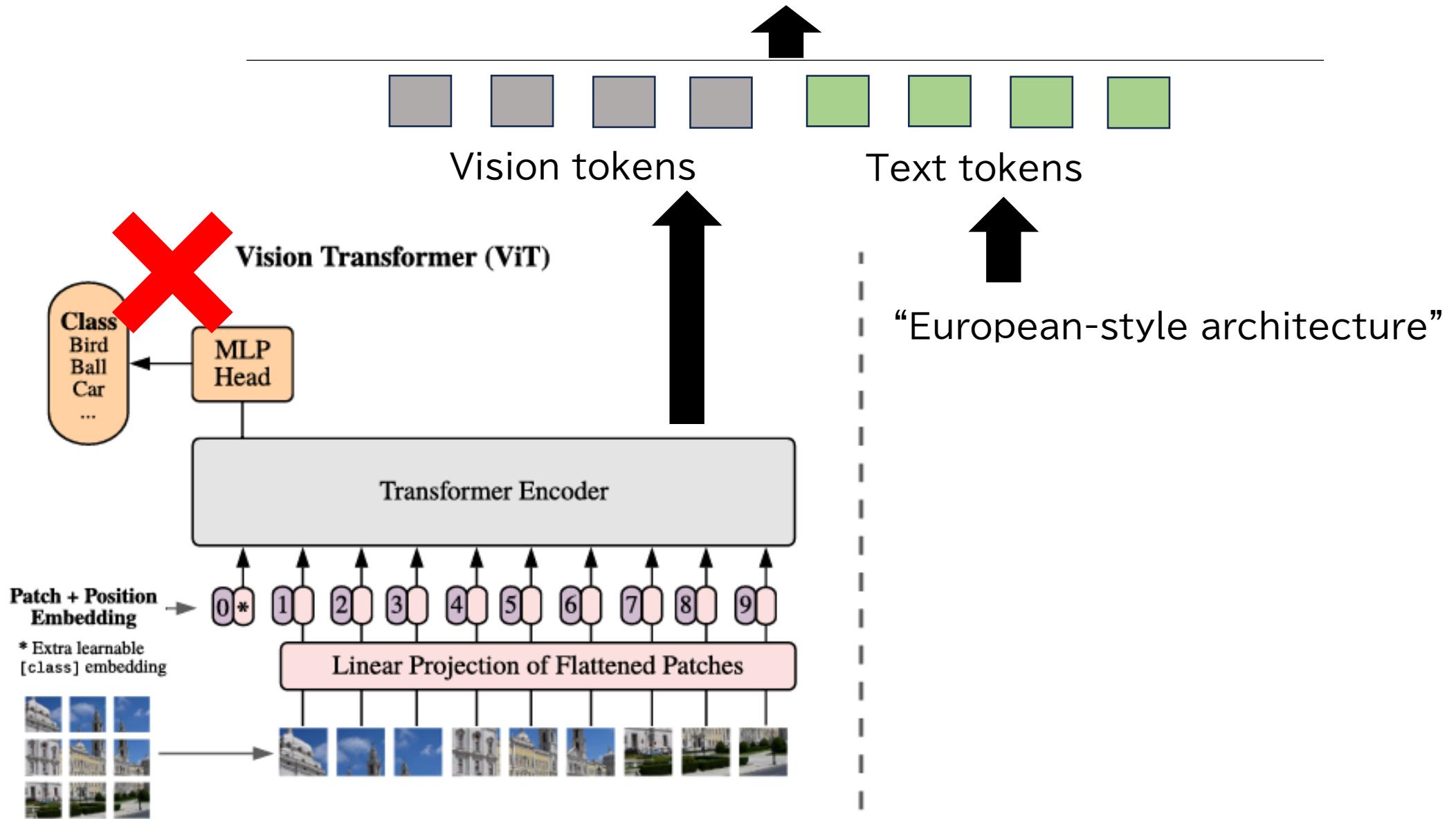


VIT

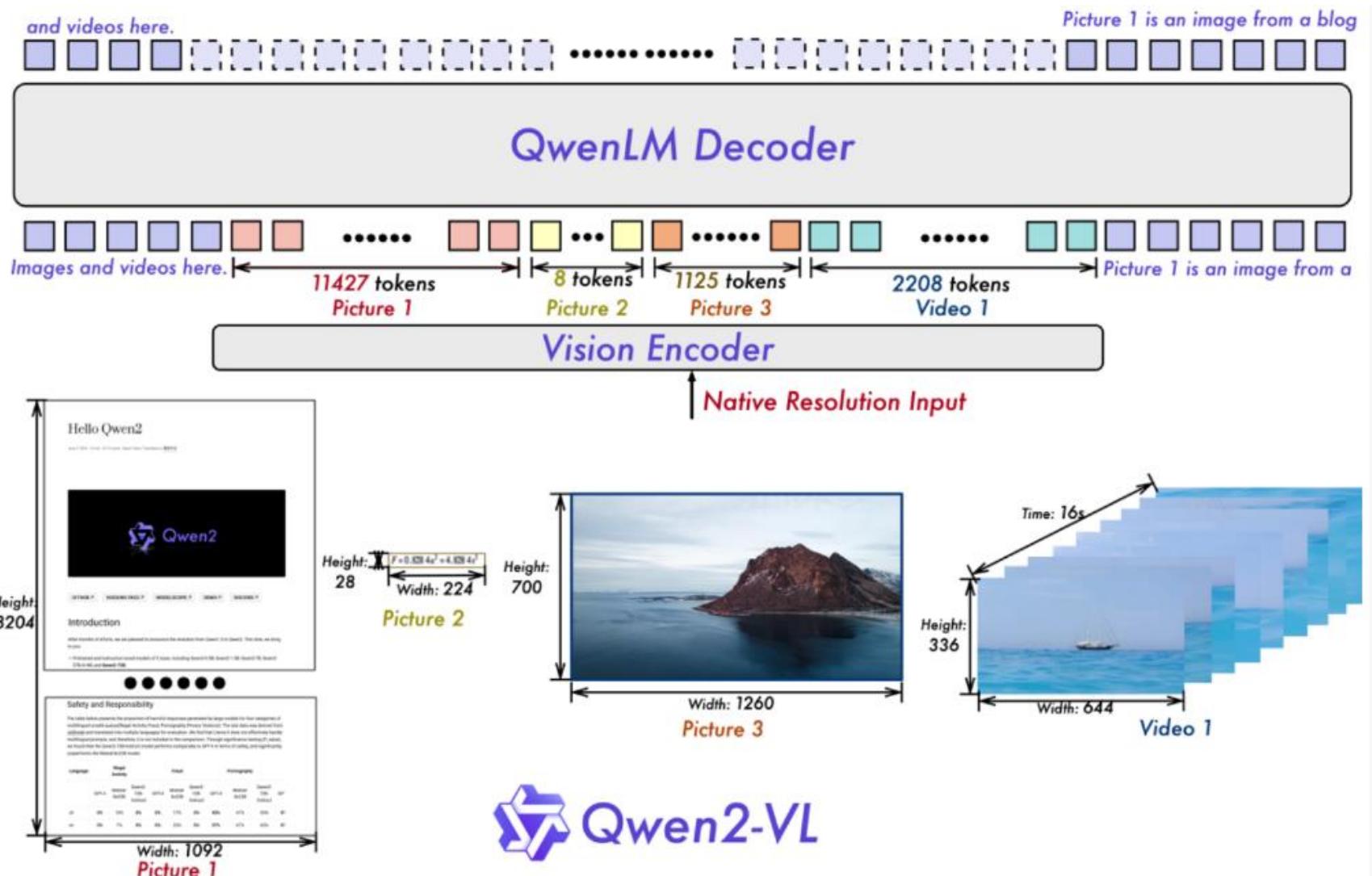


CNN

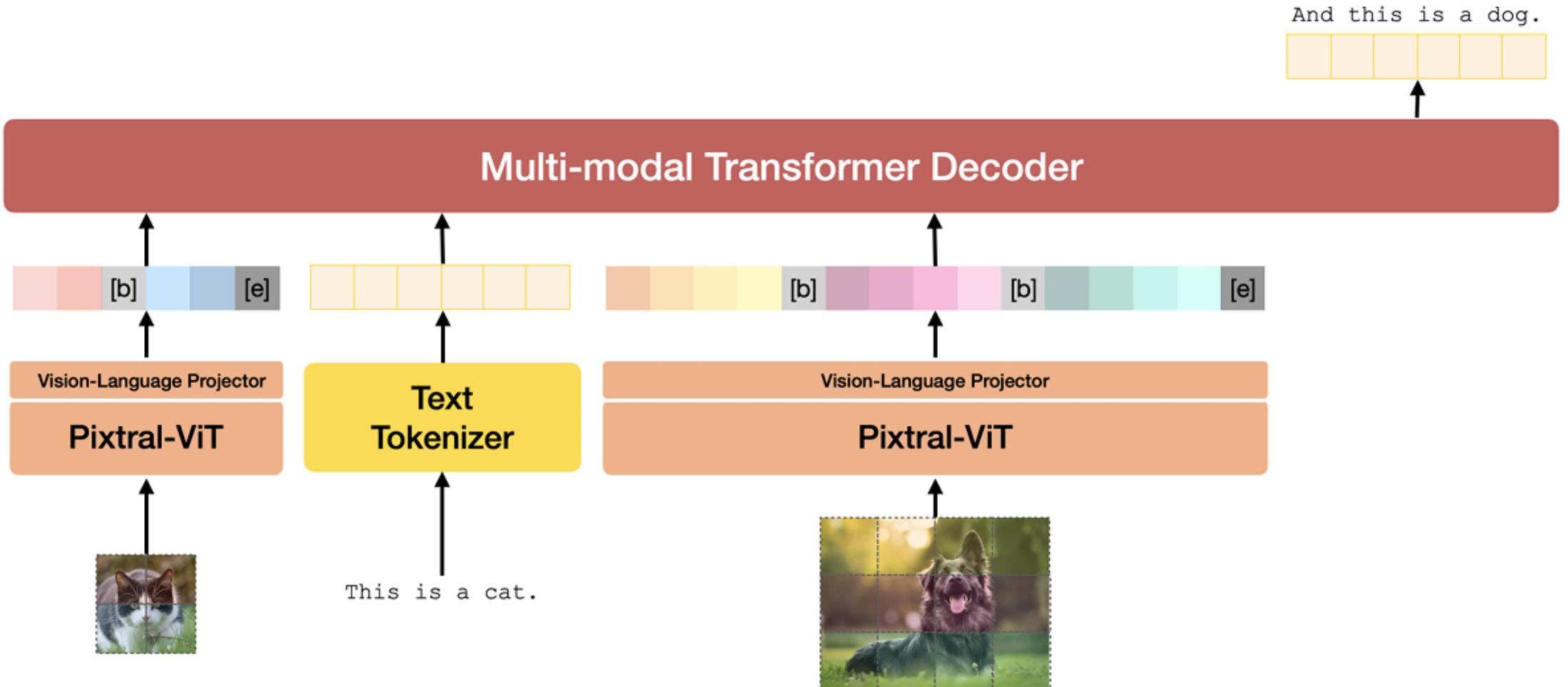
# VLM (Vision-Language Mode)



# Qwen2-VL



# Pixtral



# Motivation of VLM for table

- Many real tables
  - Can not represent as text: markdown, JSON, HTML, latex, etc…
  - Understand with the text and vision features picture, marks, symbols in the document
  - Understand table in a vision way is a natural choice

| Share-Based Compensation  |                    |                  |                   |                  |  |  |
|---|--------------------|------------------|-------------------|------------------|--|--|
| The following table shows share-based compensation expense and the related income tax benefit included in the Condensed Consolidated Statements of Operations for the three- and six-month periods ended March 30, 2024 and April 1, 2023 (in millions):  |                    |                  |                   |                  |  |  |
|   | Three Months Ended |                  | Six Months Ended  |                  |  |  |
|   | March 30,<br>2024  | April 1,<br>2023 | March 30,<br>2024 | April 1,<br>2023 |  |  |
| Share-based compensation expense  | \$ 2,961           | \$ 2,686         | \$ 5,961          | \$ 5,591         |  |  |
| Income tax benefit related to share-based compensation expense  | \$ (663)           | \$ (620)         | \$ (1,898)        | \$ (1,798)       |  |  |
| As of March 30, 2024, the total unrecognized compensation cost related to outstanding RSUs was \$24.7 billion, which the Company expects to recognize over a weighted-average period of 2.7 years.  |                    |                  |                   |                  |  |  |
| Note 9 – Contingencies  |                    |                  |                   |                  |  |  |
| The Company is subject to various legal proceedings and claims that have arisen in the ordinary course of business and that have not been fully resolved. The outcome of these contingencies is currently uncertain. In the opinion of management, there was not at least a reasonable possibility that the Company may have incurred a material loss, or a material loss greater than a recorded accrual, concerning loss contingencies for asserted legal and other claims. |                    |                  |                   |                  |  |  |
| Note 10 – Segment Information and Geographic Data   |                    |                  |                   |                  |  |  |
| The following table shows information by reportable segment for the three- and six-month periods ended March 30, 2024 and April 1, 2023 (in millions):  |                    |                  |                   |                  |  |  |
|   | Three Months Ended |                  | Six Months Ended  |                  |  |  |
|   | March 30,<br>2024  | April 1,<br>2023 | March 30,<br>2024 | April 1,<br>2023 |  |  |
| Americas:   |                    |                  |                   |                  |  |  |
| Net sales   | \$ 37,273          | \$ 37,784        | \$ 87,703         | \$ 87,062        |  |  |
| Operating income  | \$ 15,074          | \$ 13,927        | \$ 35,431         | \$ 31,791        |  |  |
| Europe:   |                    |                  |                   |                  |  |  |
| Net sales   | \$ 24,123          | \$ 23,945        | \$ 54,520         | \$ 51,626        |  |  |
| Operating income  | \$ 9,991           | \$ 9,368         | \$ 22,702         | \$ 19,385        |  |  |
| Greater China:  |                    |                  |                   |                  |  |  |
| Net sales   | \$ 16,372          | \$ 17,812        | \$ 37,191         | \$ 41,717        |  |  |
| Operating income  | \$ 6,700           | \$ 7,531         | \$ 15,322         | \$ 17,968        |  |  |
| Japan:  |                    |                  |                   |                  |  |  |
| Net sales   | \$ 6,262           | \$ 7,176         | \$ 14,029         | \$ 13,931        |  |  |
| Operating income  | \$ 3,135           | \$ 3,394         | \$ 6,954          | \$ 6,630         |  |  |
| Rest of Asia Pacific:   |                    |                  |                   |                  |  |  |
| Net sales   | \$ 6,723           | \$ 8,119         | \$ 16,885         | \$ 17,654        |  |  |
| Operating income  | \$ 2,806           | \$ 3,268         | \$ 7,385          | \$ 7,119         |  |  |

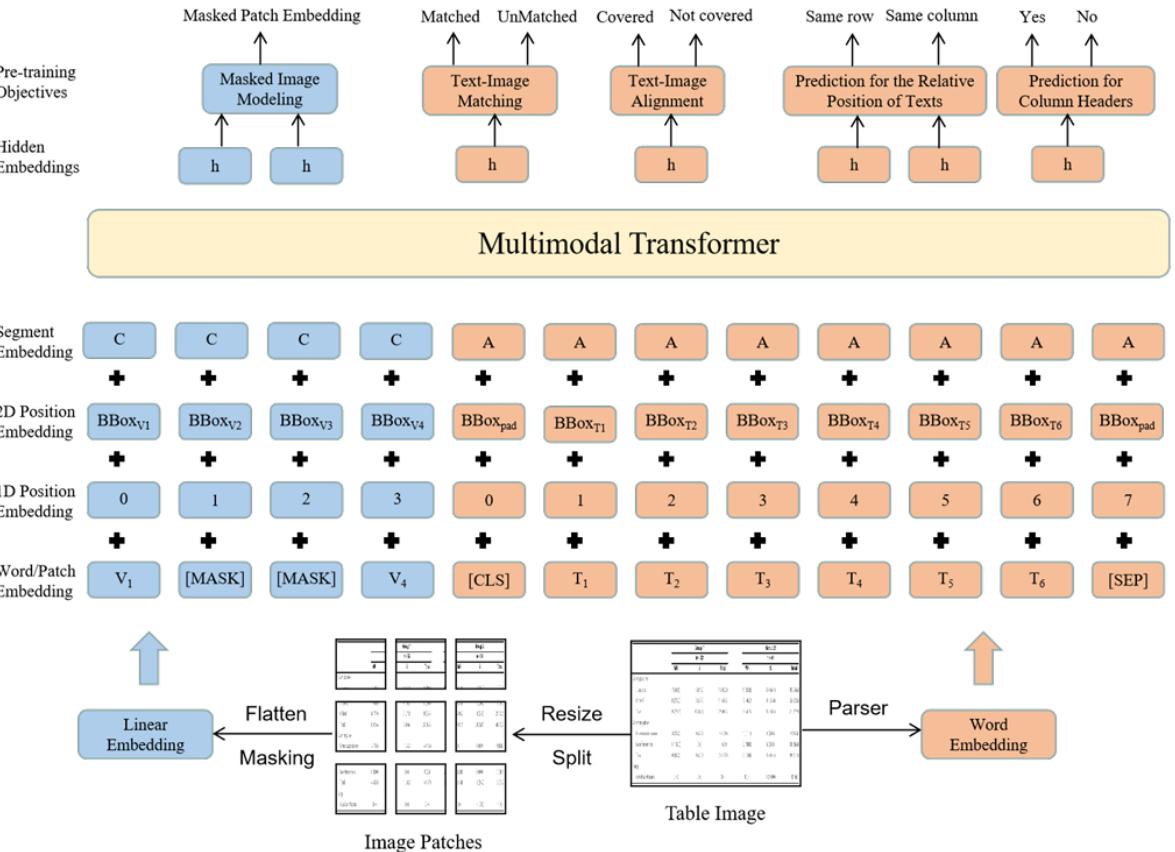
## FEATURE COMPARISON

vertex42  
© 2017 Vertex42 LLC

| FEATURES                             | PRODUCT NAME 1 | PRODUCT NAME 2 | PRODUCT NAME 3 | PRODUCT NAME 4 | PRODUCT NAME 5 | PRODUCT NAME 6 |
|--------------------------------------|----------------|----------------|----------------|----------------|----------------|----------------|
| PRICE                                | FREE           | \$25-\$30      | \$45-\$60      | \$199          | FREE           | \$500          |
| Unicode ✓ or ✗ symbols               | ✓              | ✗              | ✓              |                | ✗              | ✓              |
| Custom Icon Sets                     | ○              | ●              | ✖              | ✓              |                | ✗              |
| Up / Down Arrow Icons                | ▼              | ▬              | ▲              |                | ▼              | ▲              |
| Check / X Icons                      | ✗              | !              | ✓              | ✓              |                | ✓              |
| Flag Icons                           | 🚩              | 🚩              | 🚩              | 🚩              | 🚩              | 🚩              |
| Ratings                              | ★★★★★          | ★              | ★★             | ★★★            | ★★★★★          | ★★★★★          |
| Unicode ★ symbol                     | ★★★★★          | ★              | ★★             | ★★★            | ★★★★★          | ★★★★★          |
| Diamond ♦ symbol                     | ♦              | ♦♦             | ♦♦♦            | ♦♦             | ♦              | ♦♦♦            |
| Icon Set: Star (0,1,2)               | ☆              | ☆              | ☆              | ☆              |                | ☆              |
| Icon Set: Quarter Circle (0,1,2,3,4) | ○              | ○              | ○              | ○              | ○              | ○              |
| Icon Set: Boxes (0,1,2,3,4)          | □              | □              | □              | □              |                | □              |
| Icon Set: Bars (0,1,2,3,4)           | ▨              | ▨              | ▨              | ▨              | ▨              | ▨              |
| Numeric & Text Specifications        | 30             | 25             | 15             | 25             | 5              | 90             |
| Data Bars                            | 30             | 25             | 15             | 25             | 5              | 90             |
| Custom number formats                | 128 GB         | 64 GB          | 256 GB         | 512 GB         | 32 GB          | 1024 GB        |
| Basic Numeric Entry                  | 4.23           | 1.23           | 5              | 6.2            | 1.54           | 3.52           |
| Yes / No / na                        | No             | Yes            | No             | Yes            | na             | Yes            |
| Numbers with Different Units         | 2.3 TB         | 470 GB         | 125 GB         | 64 GB          | 512 MB         | 128 GB         |

# TableVLM(ACL23)

- Table detection and structure recognition task
  - Input: table image
  - Output: table html
- Encoder Pretrain tasks (12 layers)
  - Text-Image alignment
    - Predict if an image contain a text
  - Text-Image matching
    - Predict if image and text matched
  - Masked image modeling
- Decoder Pretrain task (4 layers)
  - Generate HTML



# TableVLM(ACL23)

- Datasets and experimental results

| Datasets            | Source   | Format     | Sizes  |
|---------------------|--|------------|--------|
| Marmot              | e-Books and Citeseer website   | bmp, xml   | 958    |
| ICDAR 2013          | European Union and US Government websites                            | pdf, xml   | 150    |
| ICDAR 2019          | modern and archival documents with various formats                   | jpg, xml   | 3.6k   |
| TableBank           | Word and Latex documents on the internet                             | jpg, HTML  | 145k   |
| SciTSR              | LaTeX source files   | pdf, Latex | 15k    |
| PubTabNet           | scientific articles in PMCOA   | png, HTML  | 568k   |
| TabLeX              | scientific paper from arXiv  | jpg, Latex | 3,00k  |
| FinTabNet           | annual reports of the S&P 500 companies                              | png, HTML  | 112k   |
| SynthTabNet         | synthetically generated based on Tablebank, PubTabNet, and FinTabNet | png, HTML  | 600k   |
| ComplexTable (ours) | synthetically generated by an auto HTML table creator                | png, HTML  | 1,000k |

Table 1: Existing public datasets available and the constructed ComplexTable dataset for table structure recognition.

| Model       | Dataset      | Simple       | Complex      | All          |
|-------------|--------------|--------------|--------------|--------------|
| WYGIWS      | TableBank    | 86.4         | --           | 86.4         |
| EDD         | TableBank    | 86.0         | --           | 86.0         |
| LGPMA       | TableBank    | 88.7         | --           | 88.7         |
| Master      | TableBank    | 89.4         | --           | 89.4         |
| TableFormer | TableBank    | 89.6         | --           | 89.6         |
| TableVLM    | TableBank    | <b>90.2</b>  | --           | <b>90.2</b>  |
| LGPMA       | PubTabNet    | 97.88        | 94.78        | 96.36        |
| Master      | PubTabNet    | 97.90        | 94.68        | 96.32        |
| TableFormer | PubTabNet    | <b>98.5</b>  | 95.0         | 96.8         |
| TableVLM    | PubTabNet    | 98.31        | <b>95.53</b> | <b>96.92</b> |
| LGPMA       | ComplexTable | 90.54        | 86.87        | 88.76        |
| Master      | ComplexTable | 92.17        | 88.79        | 90.21        |
| TableVLM    | ComplexTable | <b>94.73</b> | <b>90.43</b> | <b>92.18</b> |

Table 2: The tree-edit-distance-based similarity (TEDS) of table structure recognition on TableBank, PubTabNet and ComplexTable datasets. A table is categorized as a simple table if it lacks multi-column or multi-row cells; otherwise, it is classified as a complex table. It is worth noting that the TableBank dataset does not include any complex tables.

# PixT3 (ACL24)

- **Table to text task**
  - Input: table image
  - Output: text description of table
- Encoder-decoder Pretrain task
  - Generate row and column with a cell
- Finetune different models on specific table to text tasks

Table:

|       |     |     |       |
|-------|-----|-----|-------|
| oY    | io  | HG  | eG2S  |
| Z4iku | 01  | aRU | mubk6 |
| URa   | dAF | I   |       |
| I86   | GAe | 0b  | sUr5  |
| L1    | 3   | Vf1 | Svaq2 |

Target:

```
<<<dAF><<<URa><I>>><<<io><01><GAe>
<3>><<HG><aRU><0b><Vf1>>>>
```

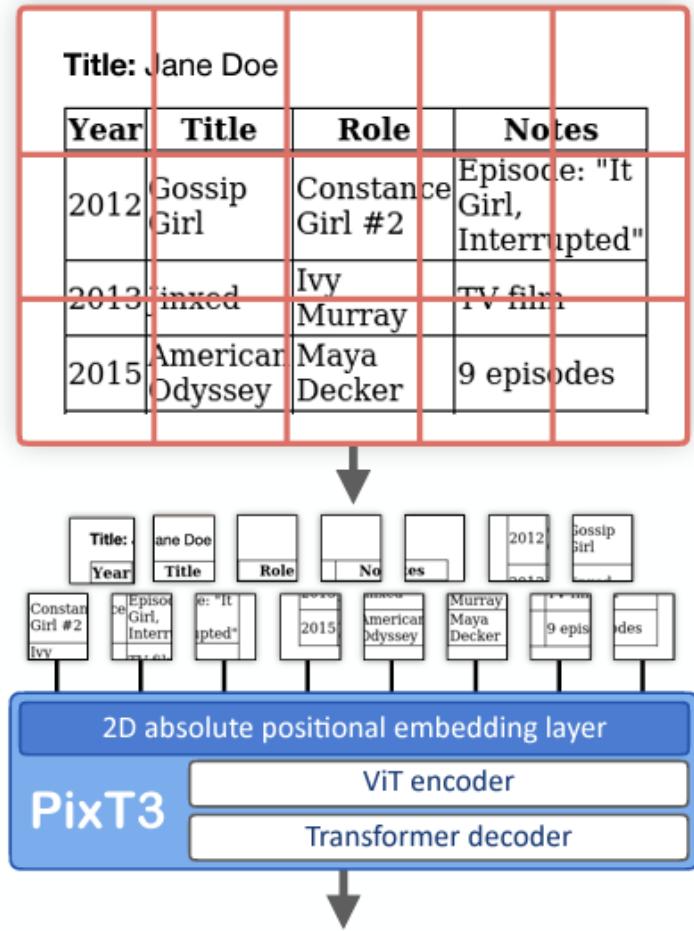
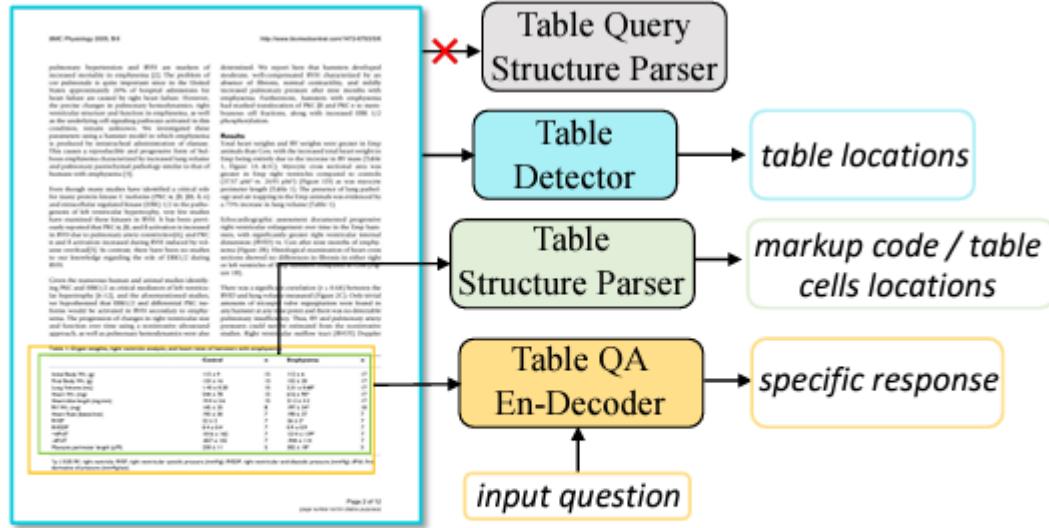


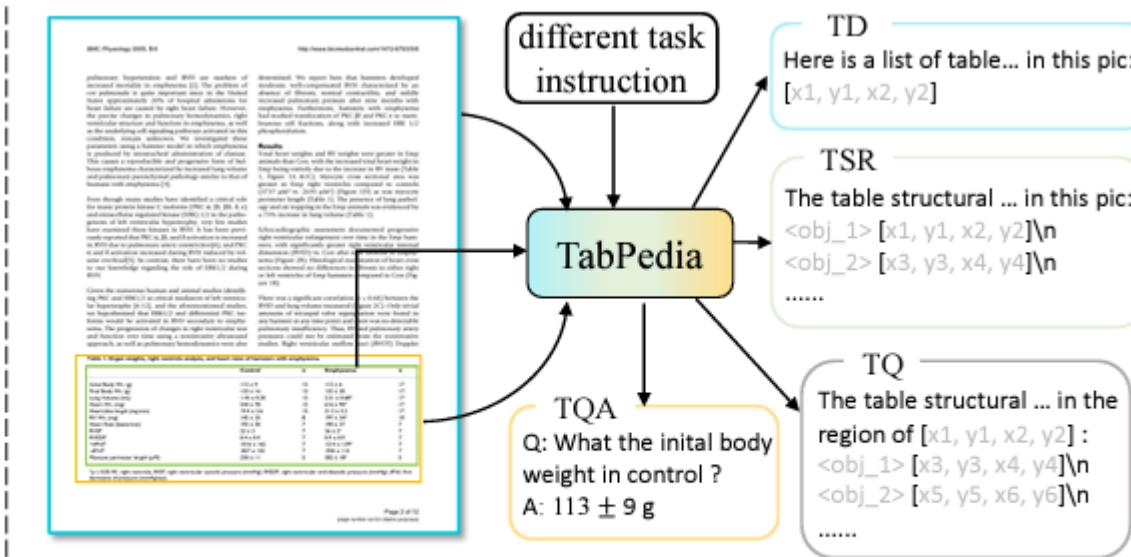
Figure 2: Overview of PixT3 generation model.

# Tabpedia (arXiv24)

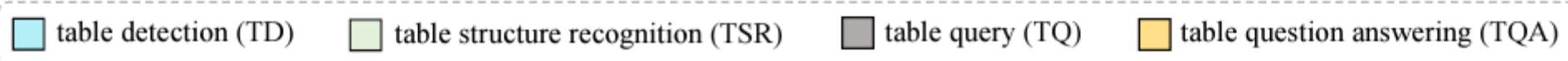
- multi-task: table detection & structure recognition, tableQA



(a) Previous task-specific pipelines

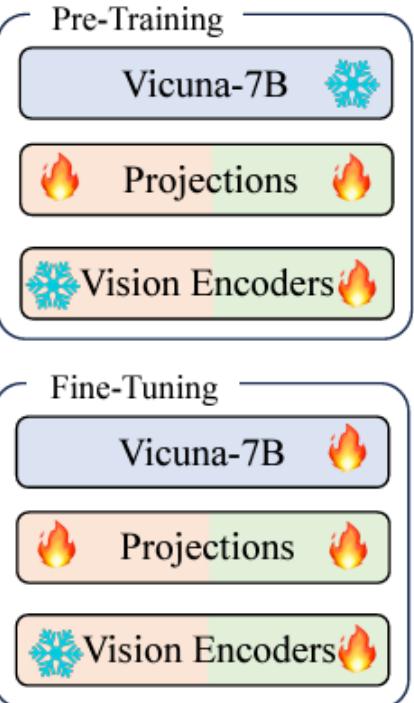
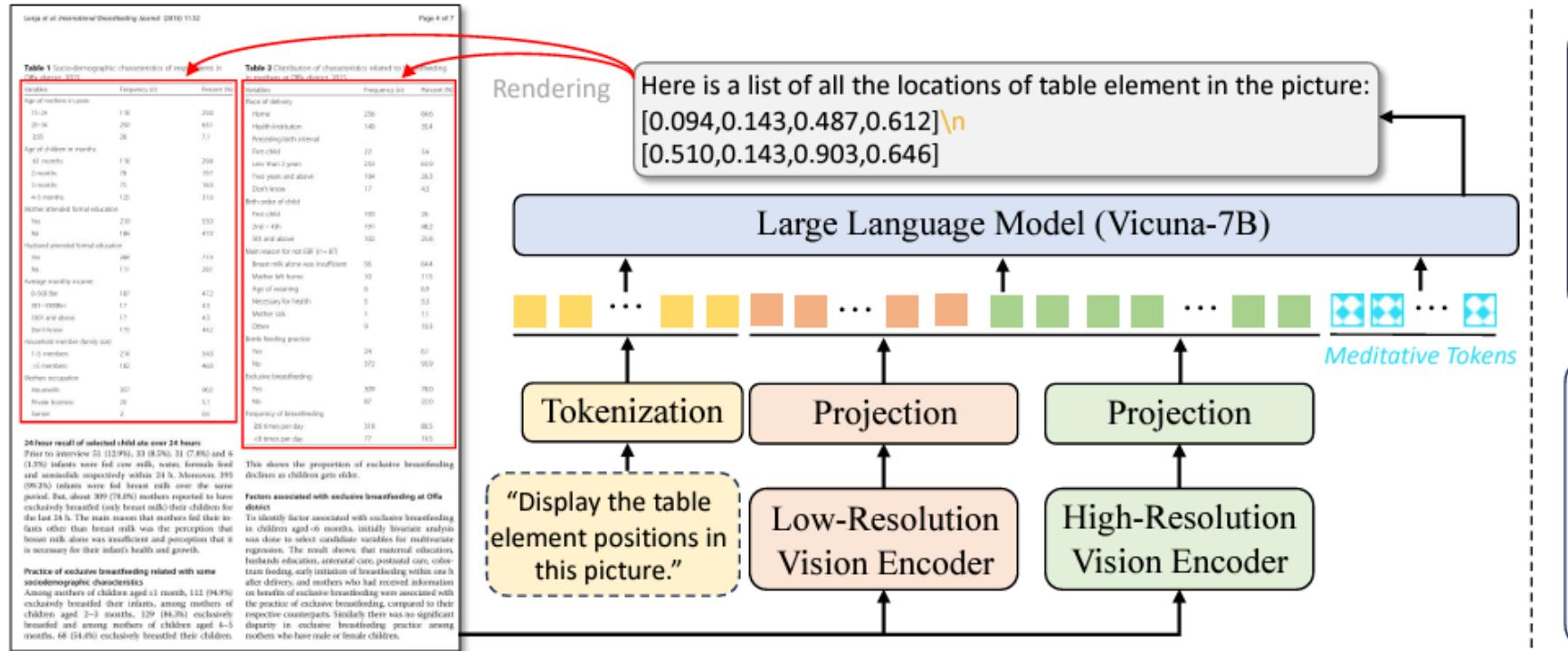


(b) Our proposed TabPedia



# Tabpedia (arXiv24)

- Architecture



# Tabpedia (arXiv24)

- Training data and tasks

Table 1: Summary of training data statistics in the fine-tuning stage.

| Dataset   | Subset       | Task    | Num  |
|-----------|--------------|---------|------|
| PubTab1M  | PubTab1M-Det | TD      | 460k |
|           | PubTab1M-Str | TSR,TQA | 759k |
|           | PubTab1M-Syn | TQ      | 381k |
| FinTabNet | –            | TSR,TQA | 78k  |
| PubTabNet | –            | TSR     | 434k |
| WTQ       | –            | TQA     | 1k   |
| TabFact   | –            | TQA     | 9k   |

Table 2: Different task types and their instruction examples.

| Task | Example   |
|------|---|
| TD   | “Give me the areas where table element’s locations in this picture.”                        |
| TSR  | “Parse the structural information of the cropped table in this picture.”                    |
| TQ   | “Parse the table structure within the region [0.095, 0.673, 0.869, 0.851] in this picture.” |
| TQA  | “What was the lowest stock price in the fourth quarter of 2010?”                            |

# Table-Llava (ACL24)

- A general VLM for multitasks

Table images of different types

| Outcome   | No. | Date             | Tournament    | Surface | Opponent in the final                   |
|-----------|-----|------------------|---------------|---------|---|
| Runner-up | 1.  | 2 September 2001 | Mostar        | Clay    | Adriana Baršić<br>Stefanie Weis         |
| Winner    | 2.  | 27 January 2002  | Courmayeur    | Hard    | Rita Degrè-Esposti                      |
| Winner    | 3.  | 17 February 2002 | Bergame       | Hard    | Dinara Safina                           |
| Runner-up | 3.  | 31 March 2002    | Rome – Parigi | Clay    | Ainika Goni-Blanco<br>Laurence Andretto |
| Runner-up | 4.  | 23 June 2002     | Gorizia       | Clay    | Sophie Lefèvre                          |
| Winner    | 4.  | 11 August 2002   | Rimini        | Clay    | Magdalena Zdeňovcová                    |
| Winner    | 5.  | 26 January 2003  | Grenoble      | Hard    | Olga Barabanshikova                     |
| Winner    | 6.  | 16 February 2003 | Southampton   | Hard    | Dinara Safina                           |
| Winner    | 7.  | 23 February 2003 | Redbridge     | Hard    | Magdalena Zdeňovcová                    |
| Winner    | 7.  | 23 March 2003    | Castellon     | Clay    | Olga Barabanshikova                     |
| Winner    | 7.  | 2 November 2003  | Poitiers      | Hard    | Roberta Vinci                           |

Web page table

| fiscal year           | all functions | national defense |
|-----------------------|---------------|------------------|
| 2010 actual           | 148962        | 86789            |
| 2011 actual           | 144379        | 83226            |
| 2012 actual           | 143737        | 79875            |
| 2013 actual           | 132477        | 70781            |
| 2014 actual           | 136159        | 70992            |
| 2015 actual           | 138544        | 72950            |
| 2016 preliminary      | 148999        | 78669            |
| 2017 proposed         | 153920        | 80480            |
| % avg growth 2010-13b | -3.8          | -6.6             |

Excel table

| year | starts | wins | top 5 | top 10 | poles | avg start | avg finish | winnings |
|------|--------|------|-------|--------|-------|-----------|------------|----------|
| 1985 | 1      | 1    | 1     | 1      | 0     | 14.0      | 1.0        | 32400    |
| 1986 | 2      | 0    | 1     | 1      | 1     | 20.0      | 21.5       | 6600     |
| 1991 | 5      | 1    | 1     | 2      | 1     | 19.2      | 18.4       | 27380    |
| 1992 | 5      | 0    | 3     | 3      | 0     | 18.6      | 12.0       | 32760    |
| 1993 | 20     | 0    | 7     | 10     | 1     | 15.6      | 15.7       | 117371   |
| 1994 | 20     | 4    | 9     | 14     | 0     | 16.6      | 11.5       | 215438   |
| 1995 | 19     | 1    | 9     | 11     | 1     | 17.6      | 12.5       | 182111   |

Markdown table

Input requests of different tasks  
(task description, required output format, task-related input)

Please provide a brief answer to the question below based on the table regarding 'Karolina Šprem'. Format your final answer as a JSON, using the structure {"answer": [<a list of answer strings>]}.

karolina sprem beat laurence andreto in 2002 but what other french player did she beat in 2003?

For the cell positioned in the 10th row and the 1st column of this table, provide its content. Output the target cell content as JSON in the form at {"row\_id": "m", "column\_id": "n", "cell\_value": "<Corresponding Cell Value>" }.

Give you an excel table about 'federal budget authority for r&d and r&d plant, by budget function category: fy 2010-17', please answer the following question ..... Question: the president's proposed budget for the federal government in fy 2017 calls for \$153.9 million in funding for r&d and r&d plant, how many million dollars of increase over the previous year?

This is a table picture. Can you figure out the row and column numbers for this particular table? Format your final answer as a JSON, using the structure {"row\_number": "m", "column\_number": "n"}.

Analyze this table picture about 'terry labonte' and classify whether the given table substantiates or conflicts with the statement. Return the result as JSON in the format {"answer": "<YOUR ANSWER>"}, e.g., {"answer": "deputed"}.

Statement: in 1995 , labonte motorsports had the most wins it had in a my single year.

This picture illustrates a table. Please represent this table with the markdown-format in text.

Output responses  
(final answer, chain-of-thoughts)

Based on the table, the final answer is {"answer": ["Sophie Lefèvre"]}.

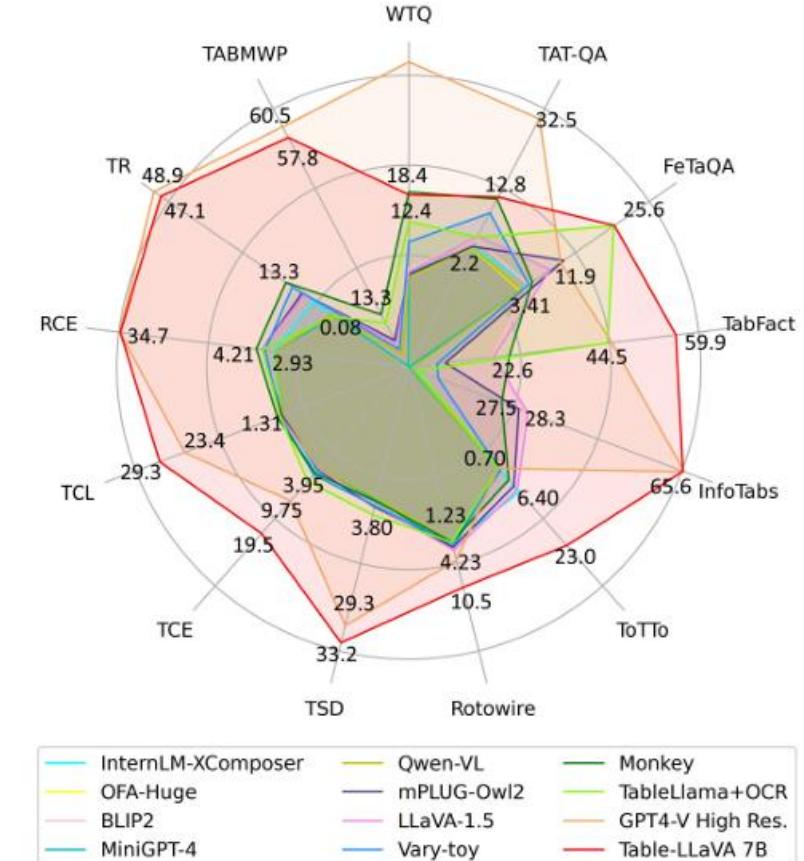
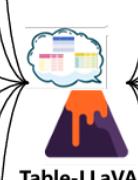
The target cell value in the 10th row and the 1st column is {"row\_id": "10", "column\_id": "1", "cell\_value": "Winner"}.

Based on the table information, we need to conduct the following computations:  
 $(153920.0 - 148999.0) = 4921$   
Therefore, the final answer is {"answer": ["4921"]}.

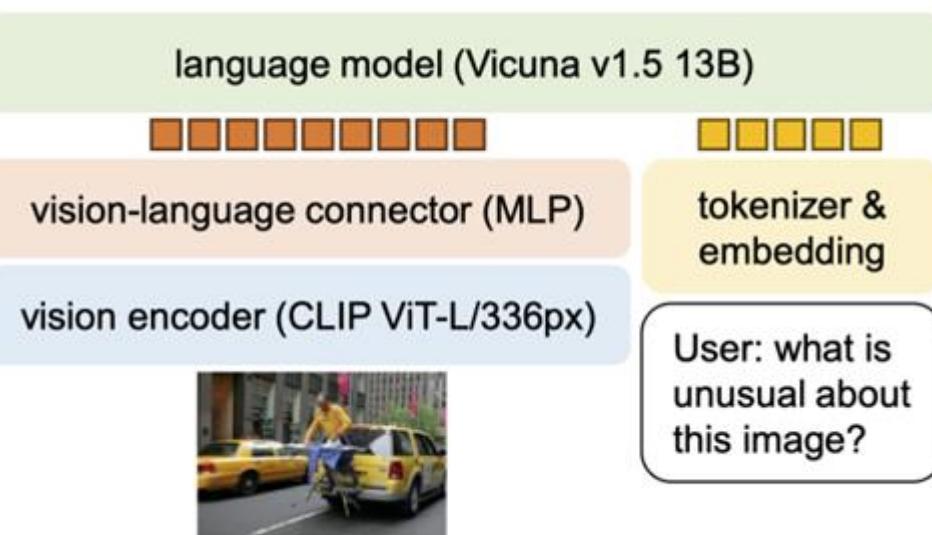
There are 11 rows and 3 columns in the table i.e. So, the final answer is {"row\_number": "11", "column\_number": "3"}.

Taking into account the table information, the table conflicts with the validity of the statement. Therefore, the final answer: {"answer": "conflicts with"}.

|   |
|---|
| year   starts   wins   top5   top10   poles |
| .....                                       |
| 1985   1   1   1   1   0   14.0   .....     |



# Table-Llava (ACL24)



Llava architecture

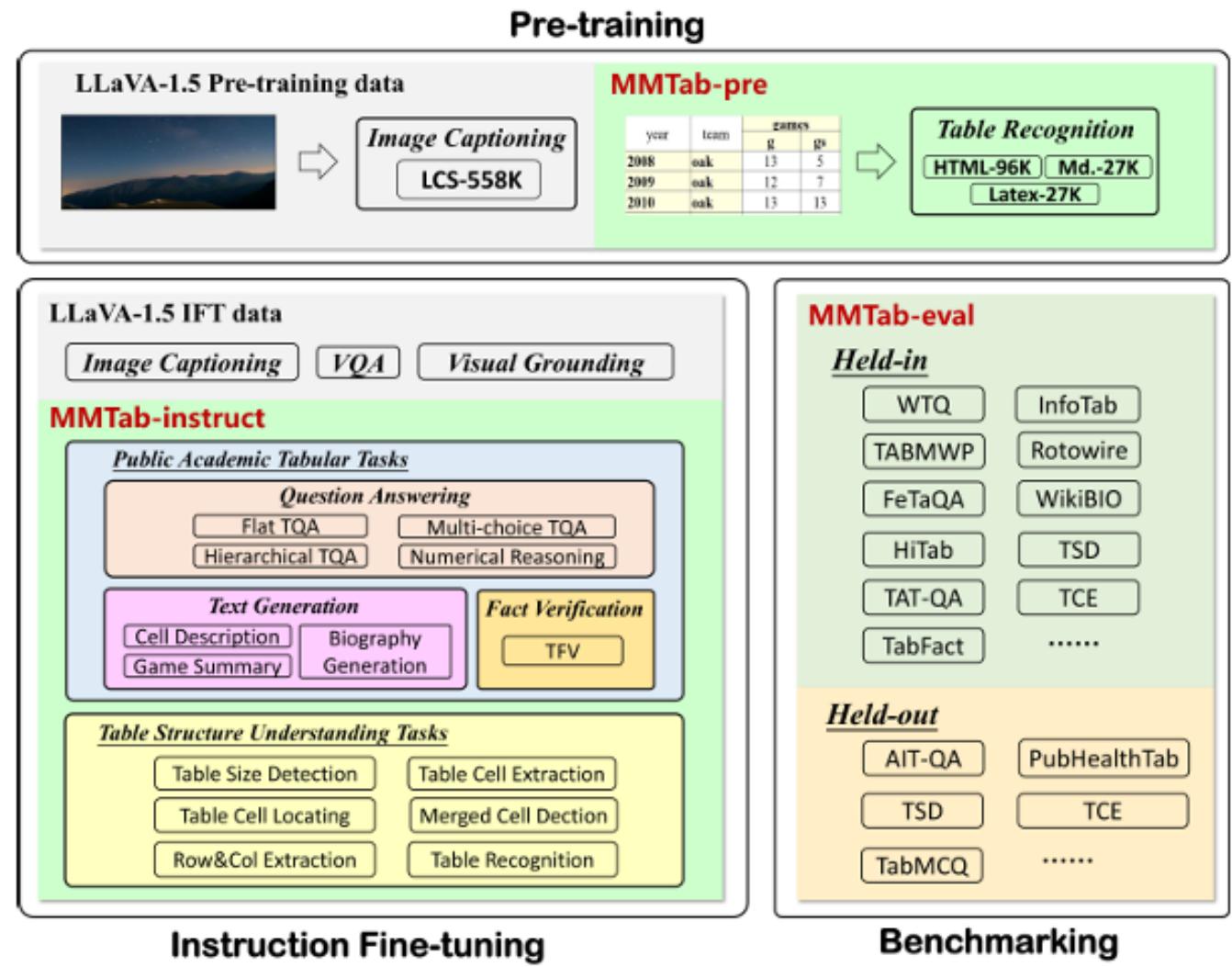


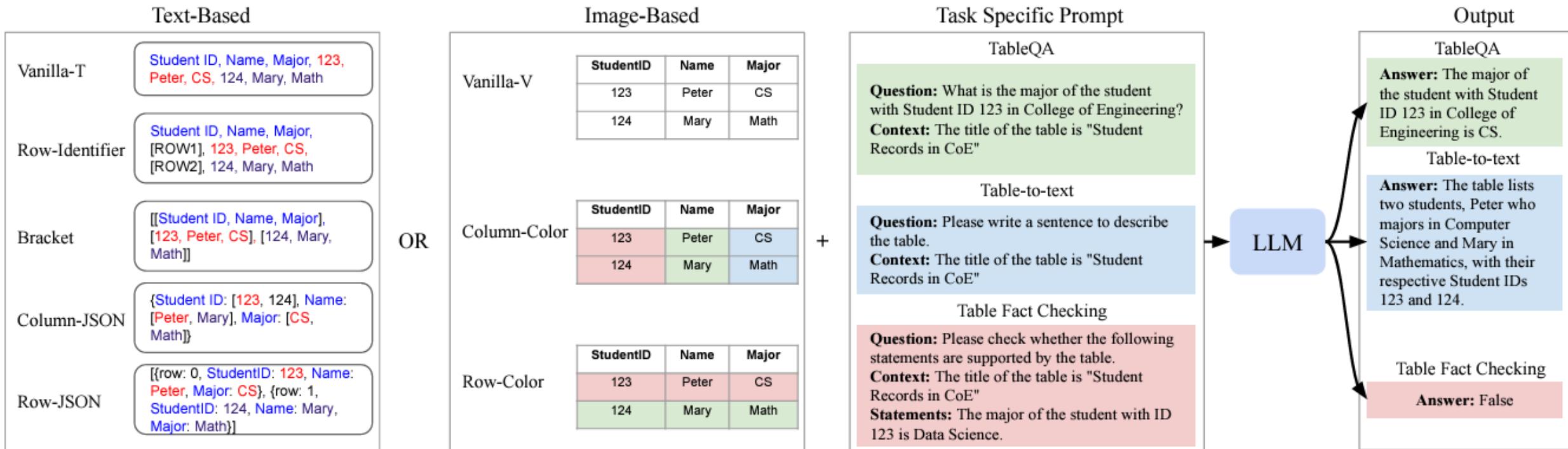
Table-Llava training

# Summary of current table VLMs

|                         | Vision encoder                                     | LM decoder backbone                    | Task  |
|-------------------------|--|--|---|
| Tableformer<br>[CVPR22] | Resnet-18(CNN) + transformer encoder (2 layers)    | Scratch transformer decoder (4 layers) | Table detection & structure recognition                 |
| TableVLM<br>[ACL23]     | ResNeXt101 (CNN) + transformer encoder (12 layers) | Scratch transformer decoder (4 layers) | Table detection & structure recognition                 |
| PixT3<br>[ACL24]        | Pix2Strct (ViT)                                    | Pix2Strct                              | Table-to-text   |
| TablePedia<br>[arxiv24] | ViT-L (ViT)<br>Swin-B (ViT)                        | Vicuna-7B                              | Table detection & structure recognition<br>TableQA      |
| Table-Llava<br>[ACL24]  | CLIP (ViT)   | Vicuna-v1.5-7B                         | Table-to-text<br>TableQA<br>Table structure recognition |

# Table as Texts or Images (ACL24)

- Evaluation the performance for different table formats.



# Table as Texts or Images (ACL24)

| Method Name    | Table Representation   |
|----------------|--|
| Vanilla-T      | $c_1, c_2, \dots, c_n, v_{(1,1)}, v_{(1,2)}, \dots, v_{(1,n)}, v_{(2,1)}, v_{(2,2)}, \dots, v_{(2,n)}, \dots, v_{(m,1)}, v_{(m,2)}, \dots, v_{(m,n)}$  |
| Row-Identifier | $c_1, c_2, \dots, c_n, [\text{ROW1}] v_{(1,1)}, v_{(1,2)}, \dots, v_{(1,n)}, [\text{ROW2}] v_{(2,1)}, v_{(2,2)}, \dots, v_{(2,n)}, \dots, [\text{ROW}m] v_{(m,1)}, v_{(m,2)}, \dots, v_{(m,n)}$  |
| Bracket        | $[[c_1, c_2, \dots, c_n], [v_{(1,1)}, v_{(1,2)}, \dots, v_{(1,n)}], [v_{(2,1)}, v_{(2,2)}, \dots, v_{(2,n)}], \dots, [v_{(m,1)}, v_{(m,2)}, \dots, v_{(m,n)}]]$ .  |
| Column-JSON    | $\{ c_1: [v_{(1,1)}, v_{(2,1)}, \dots, v_{(m,1)}], c_2: [v_{(1,2)}, v_{(2,2)}, \dots, v_{(m,2)}], \dots, c_n: [v_{(1,n)}, v_{(2,n)}, \dots, v_{(m,n)}] \}$ .   |
| Row-JSON       | $\{ \text{Row: 1, } c_1: v_{(1,1)}, c_2: v_{(1,2)}, \dots, c_n: v_{(1,n)} \}, \{ \text{Row: 2, } c_1: v_{(2,1)}, c_2: v_{(2,2)}, \dots, c_n: v_{(2,n)} \}, \dots, \{ \text{Row: m, } c_1: v_{(m,1)}, c_2: v_{(m,2)}, \dots, c_n: v_{(m,n)} \}$ . |

Text representation

| $c_1$       | $c_2$       | $\dots$  | $c_n$       |
|-------------|-------------|----------|-------------|
| $v_{(1,1)}$ | $v_{(1,2)}$ | $\dots$  | $v_{(1,n)}$ |
| $v_{(2,1)}$ | $v_{(2,2)}$ | $\dots$  | $v_{(2,n)}$ |
|             |             | $\vdots$ |             |
| $v_{(m,1)}$ | $v_{(m,2)}$ | $\dots$  | $v_{(m,n)}$ |

Vanilla-V

| $c_1$       | $c_2$       | $\dots$  | $c_n$       |
|-------------|-------------|----------|-------------|
| $v_{(1,1)}$ | $v_{(1,2)}$ | $\dots$  | $v_{(1,n)}$ |
| $v_{(2,1)}$ | $v_{(2,2)}$ | $\dots$  | $v_{(2,n)}$ |
|             |             | $\vdots$ |             |
| $v_{(m,1)}$ | $v_{(m,2)}$ | $\dots$  | $v_{(m,n)}$ |

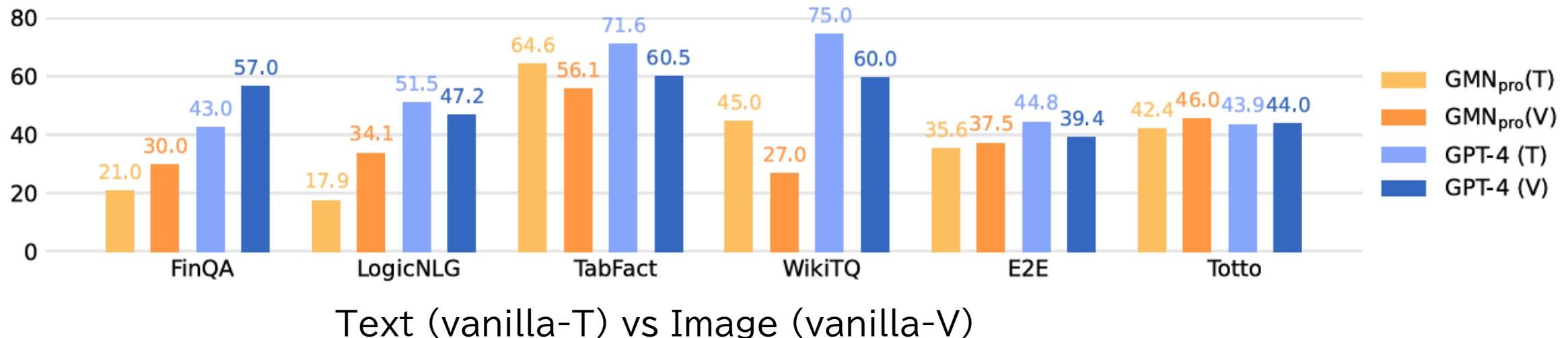
Column-Color

| $c_1$       | $c_2$       | $\dots$  | $c_n$       |
|-------------|-------------|----------|-------------|
| $v_{(1,1)}$ | $v_{(1,2)}$ | $\dots$  | $v_{(1,n)}$ |
| $v_{(2,1)}$ | $v_{(2,2)}$ | $\dots$  | $v_{(2,n)}$ |
|             |             | $\vdots$ |             |
| $v_{(m,1)}$ | $v_{(m,2)}$ | $\dots$  | $v_{(m,n)}$ |

Row-Color

image representation

# Table as Texts or Images (ACL24)



- Conclusions
  - Text is better than Image, except FinQA (math reasoning ability)
  - Best representation
    - Image: colored-row
    - Text: Bracket
  - Input both text and image for VLM not always better than either one.
  - Closed model > open model

| Datasets | T+V  | T    | V    | Metric  |
|----------|------|------|------|---------|
| WikiTQ   | 80.0 | 75.0 | 60.0 | Acc     |
| TabFact  | 64.0 | 71.6 | 60.5 |         |
| LogicNLG | 48.0 | 51.5 | 54.1 |         |
| FinQA    | 61.0 | 43.0 | 57.0 |         |
| ToTTo    | 42.4 | 43.9 | 44.0 | ROUGE-L |
| E2E      | 42.6 | 44.8 | 39.4 |         |

Table 19: GPT-4's performance when we pass the text representation (T), image representation (V) and both representation (T+V) to the model.

# TableVQA-Bench(arxiv24)

- Compare
  - VLM QA
  - VLM parse HTML + LLM QA
  - Original HTML + LLM QA
- Conclusions
  - VLM QA needs improve
  - Two phase will short the gap, but still worse than LLM QA

| Input Modality   | Model                         | VWTQ        | VWTQ-Syn    | VTabFact    | FinTabNetQA | Avg.        |
|--|-------------------------------|-------------|-------------|-------------|-------------|-------------|
| <i>Multi-modal Large Language Models (MLLMs)</i>                     |                               |             |             |             |             |             |
| Vision   | GPT-4V [1]                    | <b>42.5</b> | <b>52.0</b> | <b>68.0</b> | <b>79.6</b> | <b>54.5</b> |
|  | Gemini-ProV [23]              | 26.7        | 33.2        | 55.6        | 60.8        | 38.3        |
|  | SPHINX-MoE-1k                 | 27.2        | 33.6        | 61.6        | 36.0        | 35.5        |
|  | SPHINX-v2-1k                  | 25.3        | 28.0        | 66.8        | 31.2        | 33.7        |
|  | QWEN-VL-Chat [3]              | 19.0        | 23.2        | 60.4        | 29.6        | 28.4        |
|  | QWEN-VL [3]                   | 17.2        | 21.2        | 52.0        | 34.0        | 26.5        |
|  | SPHINX-MoE                    | 15.3        | 16.8        | 58.8        | 2.8         | 20.7        |
|  | SPHINX-v1-1k [14]             | 13.2        | 17.2        | 58.0        | 3.2         | 19.7        |
|  | mPLUG-Owl2 [26]               | 10.7        | 14.4        | 56.8        | 2.8         | 17.7        |
|  | LLaVA-1.5 [16]                | 12.4        | 12.4        | 55.6        | 0.8         | 17.7        |
|  | CogVLM-1k [24]                | 9.7         | 11.6        | 52.0        | 4.8         | 16.3        |
|  | SPHINX-v1 [14]                | 7.1         | 9.6         | 55.2        | 1.2         | 14.5        |
|  | CogAgent-VQA [7]              | 0.3         | 0.8         | 58.4        | 22.8        | 13.8        |
|  | InstructBLIP [6]              | 5.9         | 6.4         | 50.4        | 0.4         | 12.5        |
|  | BLIP-2 [13]                   | 5.2         | 5.6         | 51.6        | 0.4         | 12.2        |
|  | CogVLM [24]                   | 0.8         | 0.8         | 40.8        | 1.2         | 7.5         |
|  | CogAgent-VQA* [7]             | 37.2        | 41.2        | 58.4        | 22.8        | 39.0        |
| <i>Table Structure Reconstruction + Large Language Models (LLMs)</i> |                               |             |             |             |             |             |
| Vision   | GPT-4V [1] → GPT-4 [2]        | <b>45.2</b> | <b>55.6</b> | <b>78.0</b> | <b>95.2</b> | <b>60.7</b> |
|  | Gemini-ProV → Gemini-Pro [23] | 34.8        | 40.4        | 71.0        | 75.6        | 48.6        |
| <i>Large Language Models (LLMs)</i>                                  |                               |             |             |             |             |             |
| Text   | GPT-4 [2]                     | <b>68.1</b> | <b>69.6</b> | <b>80.0</b> | <b>98.8</b> | <b>75.5</b> |
|  | Gemini-Pro [23]               | 56.4        | 61.2        | 69.6        | 96.4        | 66.1        |
|  | GPT-3.5                       | 50.5        | 54.4        | 68.0        | 93.2        | 61.2        |
|  | Vicuna-13B [5]                | 32.8        | 39.2        | 57.6        | 84.8        | 46.7        |
|  | Vicuna-7B [5]                 | 21.5        | 34.4        | 54.0        | 68.8        | 37.0        |

# Q&A